

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Patent Application of:)
)
Hiroshi TSUDA)
) Group Art Unit: Unassigned
Serial No.: To be assigned)
) Examiner: Unassigned
Filed: January 23, 2001)



For: DOCUMENT SEARCHING APPARATUS, METHOD THEREOF, AND RECORD MEDIUM THEREOF

SUBMISSION OF CERTIFIED COPY OF PRIOR FOREIGN APPLICATION IN ACCORDANCE WITH THE REQUIREMENTS OF 37 C.F.R. §1.55

*Assistant Commissioner for Patents
Washington, D.C. 20231*

Sir:

In accordance with the provisions of 37 C.F.R. §1.55, the applicant submits herewith a certified copy of the following foreign application:

Japanese Patent Application No. 2000-028299
Filed: February 4, 2000.

It is respectfully requested that the applicant be given the benefit of the foreign filing date as evidenced by the certified papers attached hereto, in accordance with the requirements of 35 U.S.C. §119.

Respectfully submitted,

STAAS & HALSEY LLP

Date: January 23, 2001

By: _____

James D. Halsey, Jr.
Registration No. 22,729

700 Eleventh Street, N.W.
Suite 500
Washington, D.C. 20001
(202) 434-1500

JC972 U.S. PTO
09/768062



Handwritten signature

PATENT OFFICE
JAPANESE GOVERNMENT

This is to certify that the annexed is a true copy of the
following application as filed with this office.

Date of Application: February 4, 2000

Application Number: Patent Application
No. 2000-028299

Applicant(s): FUJITSU LIMITED

September 29, 2000

Commissioner,
Patent Office Kozo Oikawa

Certificate No. 2000-3079667

日 本 国 特 許 庁
PATENT OFFICE
JAPANESE GOVERNMENT

JCS972 U.S. PTO
09/768062
01/24/01

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日
Date of Application:

2000年 2月 4日

出 願 番 号
Application Number:

特願2000-028299

出 願 人
Applicant(s):

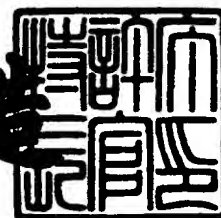
富士通株式会社

CERTIFIED COPY OF
PRIORITY DOCUMENT

2000年 9月29日

特許庁長官
Commissioner,
Patent Office

及 川 耕 造



出証番号 出証特2000-3079667

【書類名】 特許願
【整理番号】 0050081
【提出日】 平成12年 2月 4日
【あて先】 特許庁長官殿
【国際特許分類】 G06F 17/30
【発明の名称】 文書検索装置及び方法並びに記録媒体
【請求項の数】 31

【発明者】

【住所又は居所】 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

【氏名】 津田 宏

【特許出願人】

【識別番号】 000005223

【氏名又は名称】 富士通株式会社

【代理人】

【識別番号】 100074099

【住所又は居所】 東京都千代田区二番町8番地20 二番町ビル3F

【弁理士】

【氏名又は名称】 大菅 義之

【電話番号】 03-3238-0031

【選任した代理人】

【識別番号】 100067987

【住所又は居所】 神奈川県横浜市鶴見区北寺尾7-25-28-503

【弁理士】

【氏名又は名称】 久木元 彰

【電話番号】 045-573-3683

【手数料の表示】

【予納台帳番号】 012542

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9705047

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 文書検索装置及び方法並びに記録媒体

【特許請求の範囲】

【請求項 1】 リンク関係を有する文書データ群から文書データを検索する文書検索装置であって、

前記リンク関係に重みをつけて重要度であるリンク重要度を前記文書データに付与するリンク重要度付与手段と、

前記リンク重要度に基づいて前記文書データにアクセスするアクセス手段と、
を備えることを特徴とする文書検索装置。

【請求項 2】 前記リンク重要度付与手段は、前記文書データを指し示す URL (Uniform Resource Locator) 間の類似度である URL 類似度を計算する URL 類似度計算手段を備え、

前記 URL 間の類似度と文書データ間の前記リンク関係を用いて前記リンク重要度を計算する、

ことを特徴とする請求項 1 に記載の文書検索装置。

【請求項 3】 前記文書データ中からキーワードを抽出するキーワード抽出手段を更に備える、 ことを特徴とする請求項 1 または 2 に記載の文書検索装置。

【請求項 4】 前記キーワード抽出手段は、前記キーワードの前記文書データ内での出現頻度を算出し、

前記リンク重要度及び前記キーワードの出現頻度を用いて前記キーワードと前記文書データとの関連度を計算するキーワード文書関連度計算手段を更に備える

ことを特徴とする請求項 1 乃至 3 のいずれかに記載の文書検索装置。

【請求項 5】 利用者からのアクセスを監視し、アクセスログを作成する監視手段をさらに備え、

前記キーワード文書関連度計算手段は、前記キーワード出現頻度、前記リンク重要度及び前記アクセスログを用いて、関連度を計算する、

ことを特徴とする請求項 4 に記載の文書検索装置。

【請求項 6】 前記 URL 類似度並びに前記文書データのリンク数及び被リ

リンク数を用いて前記文書データの文書タイプを判別する文書タイプ判別手段をさらに備え、

前記キーワード文書関連度計算手段は、前記文書タイプに基づいて前記文書データを選択し、前記選択された文書データについて前記関連度を計算する、

ことを特徴とする請求項 4 又は 5 記載の文書検索装置。

【請求項 7】 前記抽出されたキーワードの読みや綴りから前記キーワードへアクセスできる索引を生成する索引生成手段を更に備える、

ことを特徴とする請求項 3 乃至 6 いずれか記載の文書検索装置。

【請求項 8】 利用者が前記キーワードの読みや綴りの断片を、前記索引から選択する選択手段を更に備え、

前記索引生成手段は前記キーワード文書関連度算出手段による前記関連度の高い一定個数以内の文書データを索引に入れ、

前記アクセス手段は、前記選択された前記キーワードの読みや綴りの断片に基づいて、前記文書データにアクセスする、

ことを特徴とする請求項 7 記載の文書検索装置。

【請求項 9】 ネットワークから前記文書データを収集する収集手段を更に備えることを特徴とする請求項 1 乃至 8 いずれか記載の文書検索装置。

【請求項 10】 前記リンク重要度付与手段は類似度の低い URL から多くリンクされている文書データの重要度を低くすることを特徴とする請求項 1 記載の文書検索装置。

【請求項 11】 前記リンク重要度付与手段は、重要な文書データからリンクされている、URL の類似度の低いページは重要とすることを特徴とする請求項 1 記載の文書検索装置。

【請求項 12】 前記 URL の類似度は、サーバアドレス、パス、ファイル名を含む URL の字面情報に基づいて決定される、
ことを特徴とする請求項 2 記載の文書検索装置。

【請求項 13】 対象となる文書データ集合を $DOC = \{p_1, p_2, \dots, p_N\}$

文書データ p の重要度を W_p

文書データ p のリンク先のページ集合を $R_{ef}(p)$

文書データ p のリンク元のページ集合を $R_{efed}(p)$

文書データ p と q の URL 上の類似度を $\text{sim}(p,q)$ 、相異度を $\text{diff}(p,q)=1/\text{sim}(p,q)$ とし、

文書データ p から q にリンクが張られているとした時、そのリンクの重み $lw(p,q)$

を以下で定義し、

【数 1】

$$\begin{aligned}
 lw(p,q) &= \text{diff}(p,q) / \sum_{i \in R_{ef}(p)} \text{diff}(p,i) \\
 &= \frac{1}{\text{sim}(p,q) \sum_{i \in R_{ef}(p)} \frac{1}{\text{sim}(p,i)}} \quad \dots\dots\dots(1)
 \end{aligned}$$

各ページの重要度は、各 $p \in \text{DOC}$ に対して、 C_q を定数（重要度の下限として

【数 2】

$$W_q = C_q + \sum_{p \in R_{efed}(q)} W_p * lw(p,q) \quad \dots\dots\dots(2)$$

という連立一次方程式の解として求めることを特徴とする請求項 1 記載の文書検索装置。

【請求項 1 4】 前記文書データは WWW (World Wide Web) ページであることを特徴とする請求項 1 乃至 1 3 いずれか記載の文書検索装置。

【請求項 1 5】 リンク関係を有する文書データ群の索引を生成する文書索引生成装置であって、

前記リンク関係を用いて前記文書データにリンク重要度を付与するリンク重要度付与手段と、

前記文書データ中からキーワードを抽出するキーワード抽出手段と、

前記抽出されたキーワードの読みや綴りから前記キーワードへアクセスできる索引を生成する索引生成手段と、

前記索引からキーワードの前記読みや綴りを選択された場合に、キーワード及び前記キーワードに関係する高いリンク重要度を付与された文書データにアクセスするアクセス手段と、

を備えることを特徴とする文書索引生成装置。

【請求項 1 6】 前記リンク重要度付与手段は、前記文書データを指し示す URL (Uniform Resource Locator) 間の類似度である URL 類似度を計算する URL 類似度計算手段を備え、

前記リンク重要度付与手段は、前記 URL 間の類似度と文書データ間のリンク関係を用いて前記リンク重要度を計算する、

ことを特徴とする請求項 1 5 記載の文書索引生成装置。

【請求項 1 7】 リンク関係を有する文書データ群の索引を生成する文書索引生成装置であって、

前記文書データの URL が類似しているか否かに基づいて前記文書データにリンク重要度を付与するリンク重要度付与手段と、

前記文書データ中からキーワードを抽出するキーワード抽出手段と、

前記抽出されたキーワードの読みや綴りから前記キーワードへアクセスできる索引を前記リンク重要度に基づいて生成する索引生成手段と、

を備えることを特徴とする文書索引生成装置。

【請求項 1 8】 リンク関係を有する文書データ群へのリンク集を生成するリンク集生成システムであって、

ネットワークから前記文書データを収集する収集手段と、

前記リンク関係を用いて算出された重要度であるリンク重要度を前記文書データに付与するリンク重要度付与手段と、

前記文書データ群から特定の文字列上の特徴を持つ URL を判別する URL 文

字列判別手段と、

前記リンク重要度及び前記URLの特定の文字列上の特徴に基づいて、前記文書データへのリンクを順に一定個数以内並べたリンク集を生成する索引生成手段と、

を備えることを特徴とするリンク集生成システム。

【請求項19】 前記URL類似度並びに前記文書データのリンク数及び被リンク数を用いて前記文書データの文書タイプを判別する文書タイプ判別手段をさらに備え、

前記索引生成手段は、前記文書タイプに基づいて前記文書データを選択し、前記選択された文書データについて前記リンク重要度及び前記URLの特定の文字列上の特徴に基づいて、前記リンク集を生成する、

ことを特徴とする請求項18に記載のリンク集生成システム。

【請求項20】 リンク関係を有する文書データ群から文書データを検索する文書検索方法であって、

前記リンク関係を用いて算出した重要度であるリンク重要度を前記文書データに付与する過程と、 前記リンク重要度に基づいて前記文書データにアクセスする過程と、

を含むことを特徴とする文書検索方法。

【請求項21】 前記文書データを指し示すURL (Uniform Resource Locator) 間の類似度であるURL類似度を計算する過程と、

前記URL類似度と文書間のリンク関係を用いて前記リンク重要度を計算する過程と、

を更に含むことを特徴とする請求項20に記載の文書検索方法。

【請求項22】 前記文書データ中からキーワードを抽出する過程、
を更に含むことを特徴とする請求項20又は21に記載の文書検索方法。

【請求項23】 前記キーワードの前記文書データ内での出現頻度を算出する過程と、

前記リンク重要度及び前記キーワードの出現頻度を用いて前記キーワードと前

記文書データとの関連度を計算する過程と、

を更に含むことを特徴とする請求項 2 0 乃至 2 2 いずれかに記載の文書検索方法。

【請求項 2 4】 利用者からのアクセスを監視し、アクセスログを作成する過程と、

前記キーワード出現頻度、前記リンク重要度及び前記アクセスログを用いて、関連度を計算する過程と、

を更に含むことを特徴とする請求項 2 3 に記載の文書検索方法。

【請求項 2 5】 前記 URL 類似度並びに前記文書データのリンク数及び被リンク数を用いて前記文書データの文書タイプを判別する過程と、

前記文書タイプに基づいて前記文書データを選択し、前記選択された文書データについて前記関連度を計算する過程と、

を更に含むことを特徴とする請求項 2 3 又は 2 4 に記載の文書検索方法。

【請求項 2 6】 前記抽出されたキーワードの読みや綴りから前記キーワードへアクセスできる索引を生成する過程、

を更に含むことを特徴とする請求項 2 2 乃至 2 5 いずれかに記載の文書検索方法。

【請求項 2 7】 利用者が前記キーワードの読みや綴りの断片を、前記索引から選択する過程と、

前記関連度の高い一定数以下の文書データを前記索引に入れる過程と、

前記選択された前記キーワードの読みや綴りの断片に基づいて、前記文書データにアクセスする過程と、

を更に含むことを特徴とする請求項 2 6 に記載の文書検索方法。

【請求項 2 8】 ネットワークから前記文書データを収集する過程を更に含むことを特徴とする請求項 2 0 乃至 2 7 いずれかに記載の文書検索方法。

【請求項 2 9】 リンク関係を有する文書データ群へのリンク集を自動生成する自動リンク集生成方法であって、

ネットワークから前記文書データを収集する過程と、

前記リンク関係を用いて前記文書データにリンク重要度を付与する過程と、

前記文書データ群から特定の文字列上の特徴を持つURLを判別する過程と、
前記リンク重要度及び前記URLの特定の文字列上の特徴に基づいて、前記文書データへのリンクを順に一定個数以内並べたリンク集を生成する過程と、
を含むことを特徴とするリンク集生成方法。

【請求項30】 前記URL類似度並びに前記文書データのリンク数及び被リンク数を用いて前記文書データの文書タイプを判別する過程と、

前記文書タイプに基づいて前記文書データを選択し、前記選択された文書データについて前記リンク重要度及び前記URLの特定の文字列上の特徴に基づいて、前記リンク集を生成する過程と、

を更に含むことを特徴とする請求項29に記載のリンク集生成方法。

【請求項31】 リンク関係を有する文書データ群へのリンク集を生成するコンピュータのためのプログラムを記録した記録媒体であって、

ネットワークから前記文書データを収集するステップと、

前記リンク関係を用いて前記文書データにリンク重要度を付与するステップと

、
前記文書データ群から特定の文字列上の特徴を持つURLを判別するステップと、

前記リンク重要度及び前記URLの特定の文字列上の特徴に基づいて、前記文書データへのリンクを順に一定個数以内並べたリンク集を生成するステップと、

を含む処理を前記コンピュータに実行させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、情報処理装置に蓄えられた大量の文書ファイル群を、文書の内容、文書のリンク関係及び文書の登録された場所等に基づいて検索する文書検索装置及びその方法並びに記録媒体に関する。

【0002】

【従来の技術】

今日、コンピュータネットワークの発達により、多量のオンライン文書情報があふれてきている。これらの大量のオンライン文書情報を検索及び整理するために、オンライン文書情報に索引をつけるサービスが提供されている。

【0003】

例えば、インターネットのWWW (World Wide Web) ページ検索において、ある有名なサービスでは、ディレクトリサービスを提供している。これは、階層化されたカテゴリ毎に、WWWページへのリンク集をまとめたものであり、以下の利点がある。

1. カテゴリを選択する（クリックする）だけで、利用者が閲覧を希望するWWWページへのリンクが得られる。
2. WWWページがカテゴリ毎にまとめられているため、文書検索の際に無闇に多くの情報が検索されない。
3. 人手を介してWWWページをカテゴリに登録しているので、ゴミ、つまり、利用者が閲覧を希望している情報と関係のない情報が紛れ込むことが少ない。

【0004】

これら利点のため、このサービスはインターネットにおいて非常に広く利用されている。しかし、このサービスは人手によりカテゴリの作成やWWWページの管理をしているため、運用コストがかかるという問題点がある。

【0005】

このディレクトリサービス全体の自動化にあたっては、以下の解決すべき課題がある。

1. 重要文書の選定
2. カテゴリ階層の管理、例えば、トピックの追加、削除
3. 文書からカテゴリの割付：自動分類

ここで、重要文書の選定について説明する。インターネット、イントラネット共に、WWWページは増加の一途にある。従って、似たような情報があちこちで別々に作られたりもするので、文書に含まれる文字列（キーワード）の有無や多少により検索しても結果が多すぎることになる。延いては、この多くの検索結果のうちどの情報が重要なのか分からないという問題が生じる。この重要文書の選

定に対する解決手段として以下のようなものが取られてきている。

1. 検索要求を満たす順に検索結果を並び換えする、つまり、文書中に含まれる検索キーワードの数などに基づいて、検索結果をソートし、ランキング付けをする。
2. 検索結果の視覚化によるアクセス支援を行う、つまり、検索結果の文書群を、内容の近いもの同士グルーピングする（クラスタリング）。
3. 文書の属性、例えば、サイズ、作成日時等に基づいて、検索結果をソートする。
4. その他何らかの手段で付与された文書の重要度順でソートする、例えば、リンク関係、ユーザのアクセスログ解析、第3者機関の作成した格付け等のメタデータに基づいて、検索結果をソートする。

【0006】

最後の一例として、近年WWWページのようなハイパーテキストのリンク関係を利用した文書重要度付与が、研究、サービスレベルで重要な技術となってきた。ここで、リンク関係に基づくリンク重要度付与の最も単純な実現は、「被リンク数の多い文書は重要度が高い」という直観に基づくものである。

【0007】

【発明が解決しようとする課題】

しかし、情報のナビゲーションを容易にするために、一般に同一サーバ内のWWWページは、互いにリンク関係にあることが多い。例えば、個人WWWページなどでは、「XXのトップに戻る」といった、ホームページトップへのリンクなどが多かったりする。したがって、単純な被リンク数だけでは、大量のページを抱えるサーバ（サイト）や個人が単に量が多いというだけで重要度が高いことになるという問題がある。また、検索システムが被リンク数で重要度を出していることがわかれば、意味のないページ分割や無駄なリンクだけのページを加えることで、故意に自分のページの重要度をあげることができるという問題もある。

【0008】

また、<http://www.elsevier.nl/cas/tree/store/comnet/free/www7/02/com02.htm>において閲覧可能であるWWWページにおいて、

「被リンク数の多い文書は重要度が高い。」

という直感の他に、

「重要度が高い文書からリンクされている文書は重要度が高い。」及び

「リンク先の少ないページからリンクされているページは重要度が高い。」という直感も示唆されている。

【 0 0 0 9 】

二番目の直観は、有名ディレクトリサービスで紹介される WWW ページと、名もない個人のリンク集で紹介されている WWW ページとでは、前者の方がより重要という発見に基づいたものである。三番目の直感は、例えば、1000 ものリンク先を持つリンク集からリンクされている文書よりも、50 のリンク先を持つリンク集からリンクされている文書の方が重要度が高いという考えに基づく。

これらの直観に基づいた重要度判定アルゴリズムにおいて、まず被リンク数により仮重要度を計算し、仮重要度からリンク関係により重要度を更新し、・・・という操作を収束するまで繰り返すことを行っている。

【 0 0 1 0 】

しかし、このアルゴリズムによっても、大量のページを抱えるサイトは被リンク数が多くなるため有利であり、その結果、ページの重要度を算出した場合に重要度の高いページの上位に似たようなサイトのページが並ぶという上述の問題がある。

【 0 0 1 1 】

ところで、利用者が検索をする際、検索に用いるキーワードへアクセスするインターフェースが必要となる。キーワードへのアクセスに関する従来技術として、かな漢字変換のインターフェースと関係が深い。

【 0 0 1 2 】

例えば、特開平 0 3 - 2 4 1 4 5 6 によれば、タッチパネル式のデバイスにより、かな漢字変換を実現している。この技術によれば、画面上のソフトウェアキーボードにより、キーワードの読みを全て入れてから、「変換」キーを押して漢字まじりに変換している。

【 0 0 1 3 】

また、例えば、特開平 1 0 - 1 5 4 1 4 4、特開平 1 0 - 1 5 4 0 3 3 及び <http://www.csl.sony.co.jp/person/masui/POBox/index.htm> において閲覧可能である WWW ページによれば、ペン向きテキスト入力システムを実現している。この技術によれば、画面上のソフトウェアキーボードから読みを入力するのであるが、読みの一部分を入力した時点で利用者の使用履歴に基づき随時候補を出していく。

【 0 0 1 4 】

特開平 0 3 - 2 4 1 4 5 6、特開平 1 0 - 1 5 4 1 4 4、特開平 1 0 - 1 5 4 0 3 3 及び上述の WWW ページによれば、かな漢字変換するために常に 1 文字ずつキーワードの読みや綴りを入力する必要があるため、長い文字列を入力する

また、例えば、次のような自明な読み入力インタフェースもある。この技術の例として、まず、「あ」「い」・・・といった、最初の文字毎に、キーワードリストを作り、そのキーワードリストのなかから利用者が選択する方式がある。しかし、この場合、特定の読みから始まるリスト内のキーワード数が多い場合、利用者がこのキーワードリストから特定のキーワードを選択することが困難になる。この方式の例として、銀行の自動振込機が挙げられる。

【 0 0 1 5 】

自明な読み入力インタフェース技術のさらなる例として、読みの文字を順次入力（ポインティングデバイスの場合、クリック等）して、その読みから始まるキーワード数が決まった所でかな漢字変換された文字を表示する方式もある。図 3 2 は、読みの文字を順次入力して、キーワード数が決まった所でかな漢字変換された文字を表示する方式を示す。図 3 2 では「秋葉原」を表示する場合を例として示す。図 3 2 に示すように、5 0 音の中から利用者は、キーワードの読みを順次入力する。「秋葉原」を表示させる場合、利用者は、まず「あ」を入力し、次に「き」、「は」、「ば」「ら」と順次入力する。「あきはばら」全ての音の入力が終了してキーワード数が決まったところで、かな漢字変換された文字「秋葉原」が表示される。この方式によっても、長い読みのキーワードに辿りつくまで、多くの音を入力する必要がある。

【 0 0 1 6 】



本発明は、重要文書の選定において、上述の問題、つまり、多くのページを有する特定サイトに重要度が偏るという問題を解消し、悪意を持った特定個人が重要度を操作することを困難にする文書検索装置及び方法を提供することを目的とする。

【0017】

さらに、本発明は、検索キーワードをキーワードの読みに基づいて入力する場合に、より少ない入力によりキーワードに到達でき、かつ、表示されるキーワードの候補や文書数を一定数以下にすることにより、キーワード及び文書の選択を容易にする文書検索装置及び方法を提供することを目的とする。

【0018】

さらに、本発明は、ディレクトリサービス風のインターフェースを持ち、キーワードと関連する重要文書、例えばWWWページへ迅速にアクセスできるリンク集を生成する装置及び方法を提供することを目的とする。

【0019】

【課題を解決するための手段】

本発明によれば、リンク関係を有する文書データ群から文書データを検索する文書検索装置であって、リンク関係に重みを付けた重要度であるリンク重要度を文書データに付与するリンク重要度付与手段と、リンク重要度に基づいて前記文書データにアクセスするアクセス手段を備える。

【0020】

多くのページからリンクされているページは重要であると考えられること、及び、少ないページにリンクしているページからのリンクは、多くのページにリンクしているページからのリンクよりも、重要である等と考えられることに基づいて、リンク重要度付与手段はリンク関係に重みを付けてリンク重要度を算出し、文書データに付与する。アクセス手段はこの算出されたリンク重要度に基づいて文書データにアクセスする。これにより、重要な文書データを自動的に検索することを可能とする。さらに、リンク重要度付与手段は、URL類似度計算手段を備えてもよい。URL類似度計算手段は、文書データを指し示すURL (Uniform Resource Locator) 間の類似度であるURL類似

度を計算し、リンク重要度付与手段は、URL間の類似度と文書間のリンク関係を用いてリンク重要度を計算して文書データに付与する。

【0021】

例えば同じサイト内の文書データはリンク関係を有することが多い。これら同じサイト内文書データのURLは一般にURL類似度が高い。URL類似度が高い文書データからのリンクの重みをURL類似度が低い文書データからのリンクの重みよりも低くすることにより、多量の文書データを抱えるサイトの重要度が過大に評価されることを防ぐ。延いては、より精度良く、重要な文書を検索することを可能とする。また、リンク重要度を付与する際にURL類似度を考慮するため、サイトの中の文書データのリンク数を増やしたりすることにより第三者が特定の文書の重要度を故意に上げようとするのが難しくなるという効果もある。なお、URL類似度は、サーバアドレス、パス、ファイル名を含むURLの字面情報に基づいて決定されることとしてもよい。

【0022】

また、文書検索装置は、更に文書データ中からキーワードを抽出するキーワード抽出手段を更に備えてもよい。

また、文書検索装置は、文書関連度計算手段を更に備え、上述のキーワード抽出手段は抽出されたキーワードの文書データ内での出現頻度を算出し、文書関連度計算手段はリンク重要度及びキーワードの出現頻度を用いてキーワードと文書データとの関連度を計算する。

【0023】

リンク重要度及び文書中のキーワードの出現頻度を用いて関連度を計算し、関連度の高い文書データを検索することにより、利用者が検索したいと望む文書データと内容が関連している可能性の高い、重要な文書を検索することを可能とする。

【0024】

また、文書検索装置は、利用者からのアクセスを監視し、アクセスログを作成する監視手段をさらに備え、キーワード文書関連度計算手段は、キーワード出現頻度、リンク重要度に加えてアクセスログを用いて関連度を計算することとして

もよい。関連度を算出する際にアクセスログを用いることにより、よりキーワードに関連した、より重要度の高い文書データを検索することが可能となる。

【 0 0 2 5 】

リンク重要度、キーワード出現頻度及びアクセスログを関連度の計算に用いるため、故意に特定の文書データの重要度を上げた場合であっても、そのような文書データは検索されにくくなるという効果もある。

【 0 0 2 6 】

また、文書検索装置は、URL類似度並びに文書データのリンク数及び被リンク数を用いて文書データの文書タイプを判別する文書タイプ判別手段をさらに備え、キーワード文書関連度計算手段は、文書タイプに基づいて文書データを選択し、選択された文書データについて関連度を計算することとしてもよい。

【 0 0 2 7 】

文書データには、リンクページ、コンテンツページ等、数タイプある。これらは、リンク数及び被リンク数によって判別することが可能である。この文書タイプに基づいて、あるタイプの文書データ、例えば、コンテンツページを選択し、選択された文書について関連度を計算することにより、より精度良く検索することが可能となる。

【 0 0 2 8 】

また、文書検索装置は、抽出されたキーワードの読みや綴りからキーワードへアクセスできる索引を生成する索引生成手段を更に備えてもよい。

さらに、文書検索装置は、利用者が前記キーワードの読みや綴りの断片を、前記索引から選択する選択手段を更に備えてもよい。そして、索引生成手段は索引を生成する際に、算出された関連度が高い順に、一定数以内の文書データを検索にいれる。アクセス手段は選択されたキーワードの読みや綴りの断片に基づいて、文書データにアクセスする。利用者は索引の中に含まれるデータが一定数にかぎられているために、索引から容易に選択することが可能となる。また、索引の中に含まれるデータが一定数にかぎられているため、携帯電話のような移動端末においては表示画面の広さが限られている場合に有用である。

【 0 0 2 9 】

また、文書検索装置はネットワークから文書データを収集する収集手段を更に備えてもよい。

さらに、本発明によれば、リンク関係を有するリンク集生成システムにおいて、収集手段、リンク重要度付与手段、URL文字列判別手段及び索引生成手段を備える。

【0030】

収集手段はネットワークから文書データを収集し、リンク重要度付与手段は文書データのリンク関係を用いて算出されたリンク重要度を文書データに付与し、URL文字列判別手段は文書データから特定の文字列上の特徴を持つURLを判別する。索引生成手段は、リンク重要度及びURLの特定の文字列上の特徴に基づいて、文書データへのリンクを順に一定個数以内並べたリンク集を生成する。文書データのURLの特定の文字列上の特徴は、その文書データの内容を示すことがある。例えば、J A V Aに関する文書データのURLには、J A V Aやj a v aといった文字列が用いられることがある。従って、URLの特定の文字列上の特徴は、文書データの内容を推定することに用いることができる。故に、リンク重要度及びURLの特定の文字列上の特徴に基づいて文書データへのリンク集を作成することにより、利用者が閲覧を希望する内容を含む文書データを検索することが可能なリンク集を自動的に生成することが可能となる。

【0031】

また、リンク集生成システムは、さらに文書タイプ判別手段を備えてもよい。文書タイプ判別手段は、上述のようにして文書データの文書タイプを判別し、索引生成手段は、文書タイプに基づいて文書データを選択し、選択された文書データについてリンク重要度及びURLの特定の文字列上の特徴に基づいて、リンク集を生成する。これにより、より質の高い文書データへのリンク集を生成することが可能となる。

【0032】

また、本発明の範囲は、上述の装置が実現する処理の手順からなる方法をも含む。

また、さらに、本発明の範囲は、上述の処理を前記コンピュータに実行させう

るプログラムを記録する記録媒体をも含む。

【 0 0 3 3 】

【発明の実施の形態】

以下、図面を参照しながら本発明の実施の形態を詳細に説明する。

図 1 は第 1 実施形態に係わる文書検索装置の構成を示す。図 1 の文書検索装置は、処理装置 1 1、入力装置 1 2 及び表示装置 1 3 を備える。処理装置 1 1 は、例えば、CPU (Central Processing Unit) とメモリを含み、入力装置 1 2 はキーボード、マウス等に対応し、表示装置 1 3 はディスプレイ等に対応する。

【 0 0 3 4 】

また、処理装置 1 1 は、リンク重要度付与器 2 1、キーワード抽出器 2 2、キーワード文書関連度計算器 2 3、索引生成器 2 4、索引アクセス部 2 5 及びアクセス解析器 2 6 を備える。これらは、例えば、プログラムにより記述されたソフトウェアコンポーネントに対応し、処理装置 1 1 の特定のプログラムコードセグメントに格納される。

【 0 0 3 5 】

リンク重要度付与器 2 1 は、例えば WWW ページのような文書データ 3 0 からリンク情報を抽出する。例えば、WWW ページの場合 HTML を解析し、` 富士通 トップ ` のようにアンカー (a) タグの部分を取り出す。そして、抽出されたリンク情報に基づいてリンク重要度 3 1 を算出し、算出されたリンク重要度 3 1 をキーワード文書関連度計算器 2 3 へ出力する。また、リンク重要度付与器 2 1 は、リンク元とリンク先の URL の字面の類似度である URL 類似度を算出する URL 類似度計算器 2 7 を備える。URL 類似度計算器 2 7 はリンク元とリンク先の URL 類似度を算出し、リンク重要度付与器 2 1 は、抽出されたリンク関係及び URL 類似度に基づいてリンク重要度 3 1 を算出する。

【 0 0 3 6 】

キーワード抽出器 2 2 は文書データ 3 0 からキーワードを抽出し、その結果をページキーワード 3 2 として出力する。また、キーワード抽出器 2 2 は抽出され

たキーワードが文書データ30中で出現する全出現頻度を集計することとしてよい。キーワード抽出器22は、例えば、文書データ30が日本語で記述されている場合に、形態素解析（単語切り）を行い、名詞（列）をキーワードとして抽出する。また、簡単な表記の揺れ（「コンピュータ」と「コンピューター」など）はルールや小規模なルールで統一しておく。同義語の情報は例えば、外付けで辞書などにより与えられる。

【0037】

キーワード文書関連度計算器23は、リンク重要度31、ページキーワード32及び後述のアクセスログ34に基づいて、キーワードと文書の関連度であるキーワード文書関連度を計算し、計算結果を索引生成器24へ出力する。

【0038】

索引生成器24は、キーワード文書関連度に基づいて索引データ33を生成し、生成された索引データ33を索引アクセス部25に出力する。索引データ33は例えばハイパーテキストを用いて生成される。

【0039】

索引アクセス部25は、入力装置12から入力される利用者の指示に従って、索引データ33の内容を表示装置13に表示し、利用者のアクセス状況を示す情報をアクセス解析器26に出力する。

【0040】

アクセス解析器26は、利用者のアクセス状況を示す情報を解析し、一定期間内に各キーワードからアクセスされた文書を集計したアクセスログ34を作成し、作成されたアクセスログ34をキーワード文書関連度計算器23に出力する。

【0041】

次に、図2から図5を参照しながら各種データのデータ構造について説明する。

図2は文書情報を格納するテーブル集を示す。文書情報を格納するテーブル集は、文書情報テーブル41及び被参照文書テーブル42を含む。文書情報テーブル41は、文書ID、URL、タイトル、被参照文書テーブル42へのリンク及び文書データのリンク重要度からなる。文書IDは文書データにユニークにつけ

られた数値、URLはインターネットで文書データを指す情報、タイトルはその文書データのタイトル、被参照文書テーブル42はその文書データにリンクしているリンク元の文書データの集合を格納する。被参照文書テーブル42は文書ID及びリンク先とリンク元の文書データのURL上の字面の類似度であるURL類似度の対情報を格納する。1文書毎に高々1つの被参照文書テーブルがある。文書情報テーブル41及び被参照文書テーブル42は、リンク重要度付与器21が作成するリンク重要度31に相当する。

【0042】

図3はキーワード情報を格納するテーブル集を示す。キーワード情報を格納するテーブル集はキーワードテーブル51、キーワード対応テーブル52及び出現文書テーブル53を含む。キーワードテーブル51は、キーワードID、代表語、出現文書テーブル53へのリンクからなる。代表語とは、同じキーワードIDを持つキーワードが複数ある場合に、どれを代表とするかを示すの情報である。キーワード対応テーブル52は、実際のキーワードと、読み（又はスペル）、キーワードIDとの対情報からなる。ここで、同一概念を表すキーワード（例えば、「コンピュータ」「computer」「計算機」など）には同一のキーワードID（kwID）がふられている。また、日本語のキーワードの読みについては、読みの中の「ー」を除去したり、「ぁ」や「ぃ」のような拗撥音を「あ」や「い」として表記を統一しておく。英語のキーワードについては、大文字に統一しておく。これにより、「コンピューター」や「コンピュータ」といった表記のゆれによって同じ概念を示すキーワードが異なるように扱われることを防ぎ、生成される索引においてキーワードの統一を行うことができる。出現文書テーブル53は、当該キーワードを含む文書データの文書IDとそのキーワードの出現頻度からなる。キーワード対応テーブル52及びキーワードテーブルの代表語は予め処理装置11に備えられている（不図示）。出現文書テーブル53はキーワード抽出器22により作成されるページキーワード32に相当する。

【0043】

図4は索引情報を格納するテーブル集を示す。索引情報を格納するテーブル集は、索引情報テーブル61、関連文書テーブル62及び関連キーワードテーブル

63を含む。索引情報テーブル61は、索引文字列、継続する文字列の並び、関連キーワード列からなる。索引情報テーブル61は、図3のキーワード対応テーブル52に格納されたキーワードとその読み（又はスペル）及びキーワード文書に基づいて、索引データ生成部24が後述の方法により文字列グラフを生成し、それを縮退することにより作成される。図4では、例えば、トップが「あ」「い」……。 「あ」を選ぶと次に「あいぼ」，「あお」，……。 が現れることを示している。また、文字列「あいぼ」から考えられるキーワードは「相棒」「アイボリー」であることを示している。ここでキーワードは、図4のキーワード対応テーブル52に含まれるキーワードである。関連文書テーブル62は、キーワードIDから関連する文書のIDである関連文書IDを得るためのテーブルである。キーワード文書関連度計算器23は文書関連度を計算し、その計算結果に基づいて、文書関連度の高いものから順に一定個数以内の関連文書ID列を関連文書テーブル62に格納する。関連キーワードテーブル63は、文書IDから関連キーワードIDを得るためのテーブルである。関連文書テーブル62と関連キーワードテーブル63の情報は同じで、それを転置したものである。なお、関連文書IDの詳細情報は、図2の文書情報テーブル41に格納されている。

【0044】

図5は、アクセスログ情報を格納するテーブルであるアクセスログを示す。アクセスログ71は、利用者がキーワード情報画面（後述）から文書データを選んだ際に関する情報、つまり、アクセス日時、キーワードID及び選ばれた文書データの文書IDの対情報等を格納する。アクセスログ71は、アクセス解析器26が作成するアクセスログ34に相当する。一定期間内のログを集計することにより、特定キーワードから特定文書がどのくらいアクセスされたかを得ることができる。

【0045】

以下、文書検索装置全体の動作について図6を用いて説明する。図6は索引生成処理の手順を示す。

まず、リンク重要度付与器21は、文書データからリンク情報及びURL等を抽出する。続いて、抽出された情報を文書情報テーブル41の文書ID、URL

フィールドに格納し、被参照文書テーブル 4 2 へのリンク（ポインタ）及び被参照文書テーブル 4 2 を生成する（ステップ S 1）。

【 0 0 4 6 】

リンク重要度付与器 2 1 に備えられた URL 類似度計算器 2 7 は、抽出されたリンク情報及び URL に基づいて、リンク元とリンク先の URL 類似度を計算し、被参照文書テーブル 4 2 のフィールドに格納する。

【 0 0 4 7 】

次に、リンク重要度付与器 2 1 は抽出されたリンク情報及び URL 類似度に基づいてリンク重要度を算出し、文書情報テーブル 4 1 のフィールドに格納する（ステップ S 2）。URL 類似度及びリンク重要度の算出については後述する。

【 0 0 4 8 】

キーワード抽出器 2 2 は、文書データからキーワードを抽出し、キーワード対応テーブル 5 2 のキーワードとキーワード ID のフィールド、キーワードテーブル 5 1 の全フィールド及び出現文書テーブル 5 3 の文書 ID と頻度のフィールドに格納する（ステップ S 3）。キーワードの抽出は、例えば、文書データ 3 0 が日本語で記述されている場合に、キーワード抽出器 2 2 は形態素解析（単語切り）を行い、単語切りによって得られた名詞（列）から行う。また、簡単な表記の揺れ（「コンピュータ」と「コンピューター」など）はルールや小規模な辞書で統一しておく。同義語の情報は外付けで辞書などにより与えられる。

【 0 0 4 9 】

次にキーワード抽出器 2 2 は、抽出されたキーワードの読み付けを上述のようにして統一された表記法則に基づいて行い、その結果をキーワード対応テーブル 5 2 の読みフィールド（又はスペル）に格納する（ステップ S 4）。キーワード対応テーブル 5 2 ではキーワードの表記が統一されて格納されているため、これにより、生成される索引のキーワードを統一することができる。

【 0 0 5 0 】

次にキーワード抽出器 2 2 は、抽出されたキーワードが文書データ中で出現する全出現頻度を集計し、その結果に基づいてキーワードテーブル 5 1 の出現文書のポインタを生成し、出現文書テーブル 5 3 の文書 ID 及び頻度のフィールドに

集計された頻度を格納する（ステップS5）。さらに、キーワード抽出器22は、キーワードIDの全出現頻度を集計し、上位一定個数（例えば10,000語）のキーワードを索引の対象キーワードとして決定し、当該キーワードID以外のエントリを、キーワードテーブル51及びキーワード対応テーブル52から除く。

【0051】

次に、キーワード文書関連度計算器23は、文書情報テーブル41のリンク重要度、被参照文書テーブル42のURL類似度及びアクセスログ71に基づいてキーワードと文書の関連度を示すキーワード文書関連度を計算し、キーワード文書関連度が上位から一定個数にある文書を関連文書として決定し、その決定に基づいて関連文書テーブル62及び関連キーワードテーブル63の関連文書ID列のフィールドに情報を格納する（ステップS6）。

【0052】

次に、索引生成器24は、キーワード対応テーブル52のエントリキーワードとその読み（又はスペル）に基づいて、文字列グラフを生成し、それを縮退し、その結果を索引情報テーブル61に格納する（ステップS7）。縮退については後述する。

【0053】

次に、索引生成器24は、索引情報テーブル、関連文書テーブル、関連キーワードテーブル及び文書情報テーブルに基づいて索引を作成する（ステップS8）。索引は例えばハイパーテキストで作成される。作成された索引は、表示装置13に表示されてもよい。

【0054】

生成された索引は索引アクセス部25を介して表示装置13に出力され、利用者は表示された索引を用いて入力をする。索引アクセス部25は利用者のアクセス状況を示す情報をアクセス解析器26に出力する。アクセス解析器26はアクセス状況を示す情報を解析してアクセスログ34を作成する（不図示）。

【0055】

以下、文書検索装置のリンク重要度付与器がリンク重要度を算出する手順について説明する。

まず、本実施形態において、リンク重要度付与器 21 はリンク重要度を付与する際に、その文書データのリンク関係、URL 及びキーワードを利用する。なお、リンク関係に基づいて判定される文書データの重要度をリンク重要度という。リンク重要度を判定する際の基本的な考え方は以下の通りである。

1. 類似度の低い URL から多くリンクされている文書データ（ページ）は重要である。

【0056】

例えば、一般に、同一サイト内に設けられた複数の WWW ページ（ページ）はそのサイト内の他のページにリンクされているが、それらのページの URL は相互に類似する。従って、類似度の高い URL からリンクされているページの重要度は低いと推定できるからである。

2. 多くのページからリンクされているページほど重要なページであり、重要なページからリンクされている、URL の類似度の低いページは重要である。

【0057】

例えば、有名なディレクトリサービス等及び官公庁等は多くのページからリンクされているが、このような重要なページからリンクされている文書の方が、個人が開設するページやそのコンテンツのエントリページからリンクされている文書よりも重要度が高いと考えられるからである。また、多くのページやミラーサイトを抱えるサービス（サイト）に設けられたページ等はそのサイト内でリンクされていることが多い。そのため、従来技術において、同じサイトのページが多く検索されてしまうという問題があった。しかし、そのサイト内のページの URL は例えばドメインが同じ等大抵類似しているため、「URL の類似度の低いページは重要である」という考え方を導入すれば、この問題を解消することが可能となる。

3. URL の類似度は、サーバアドレス、パス、ファイル名の全てが異なるものが最も小さく、ミラーサイトや同一サーバ内のページは類似度が高くなるように、URL の字面情報から定義する。

【0058】

上述の 3 つの考え方を導入することにより、全てのリンク関係を同等に扱わな

いでリンク重要度に応じた重みをリンク関係に与えることとしている。より具体的には、リンクの重みをリンク元とリンク先文書のURLの類似度の逆数として与えることとしている。これにより、単純にリンクされている数だけで文書の重要度を判定しすべてのリンクを同等に扱っている従来技術における問題、つまり、大量のページやミラーサイトを抱えるサーバ（サイト）や個人が単に量が多いというだけで重要度が高いことになるという問題を解決することが可能となる。また、故意にサイト内のページを増やしてリンクを設定することによりページの重要度を上げようとしても、同じサイトのページのURL類似度は高いために、ページの重要度を上げることは従来よりも困難になるという効果もある。

【 0 0 5 9 】

以下、リンク重要度付与器 2 1 におけるリンク重要度の算出についてより詳しく説明する。

リンク重要度の算出対象となるページ集合を $DOC = \{p_1, p_2, \dots, p_N\}$ 、
 ページ p のリンク重要度を W_p 、
 ページ p のリンク先のページ集合を $Ref(p)$ 、
 ページ p のリンク元のページ集合を $Refed(p)$ 、
 ページ p と q のURL類似度を $sim(p, q)$ 、
 相異度を $diff(p, q) = 1/sim(p, q)$ とすると、
 ページ p から q にリンクが張られているとした時、そのリンクの重み $lw(p, q)$ を
 以下の (1) 式で定義する。

【 0 0 6 0 】

【数 3】

$$\begin{aligned}
 lw(p, q) &= diff(p, q) / \sum_{i \in Ref(p)} diff(p, i) \\
 &= \frac{1}{sim(p, q) \sum_{i \in Ref(p)} \frac{1}{sim(p, i)}} \dots\dots\dots (1)
 \end{aligned}$$

【 0 0 6 1 】

この (1) 式から分かるように、 $lw(p,q)$ は、 p と q の URL の類似度 $sim(p,q)$ が低いほど、また、 p からのリンク数がより少ないほど大きくなる。

各ページのリンク重要度は、各 $p \in DOC$ に対して、 C_q を定数 (重要度の下限であり、ページによって異なる値を与えてもよい。) として、

【 0 0 6 2 】

【 数 4 】

$$W_q = C_q + \sum_{p \in Refed(q)} W_p * lw(p,q) \quad \dots\dots\dots (2)$$

【 0 0 6 3 】

という連立一次方程式の解として定義する。リンク重要度付与器 2 1 は、この連立一次方程式を解くことにより、リンク重要度を各ページに付与する。なお、このような連立一次方程式の解法については、既存のアルゴリズムが多数存在するため、説明は省略する。(1) 式中の URL 類似度 $sim(p,q)$ はリンク重要度付与器 2 1 に備えられた URL 類似度計算器 2 7 により、計算される (後述) 。(1) 式及び (2) 式から、上述の考え方が実現されていることを読み取ることができる。すなわち、(1) 式から類似度が低ければ、重み lw は大となるから (2) 式から類似度の低い URL から多くリンクされている文書ページは重要となる。また、(2) 式から多くのページからリンクされているページほど重要なページとなる。

【 0 0 6 4 】

さらに、(2) 式から重要なページ (W_q) からリンクされている URL の類似度の低いページ (リンクの重み lw の高い) は重要である。以下、図 7 及び図 8 を用いて、(1) 式及び (2) 式の示す考え方についてより詳しく説明する。

【 0 0 6 5 】

図 7 は、(1) 式及び (2) 式の示す内容を模式的に示す。図 7 において、まるは各ページ、矢印はリンク関係、矢印の方向はリンクの方向、矢印の太さがリ

リンクの重みを示す。図7に示すように、ページ p_1 、 p_2 及び p_3 からページ q へリンクがはられている。ページ p_1 はページ q 以外の2つのページ r_1 及び r_2 にもリンクし、ページ p_3 は他の2つのページ s_1 及び s_2 からリンクされている。

【0066】

ここで、各ページのURL類似度は、

$$\text{sim}(p_1, q) = \text{sim}(p_1, r_1) = \text{sim}(p_2, r_1) = 1$$

$$\text{sim}(p_2, q) = 2, \quad (\text{つまり、ページ } p_2 \text{ 及び } q \text{ のURLは若干類似する})$$

$$\text{sim}(p_3, q) = 1, \quad \text{sim}(s_1, p_3) = \text{sim}(s_2, p_3) = 3$$

(つまり、ページ s_1 、 s_2 及び p_3 のURLは類似する)

とする。

【0067】

(1) 式及び(2) 式を図7のような場合に適用すると、

各ページ p_1 、 p_2 、 p_3 、 s_1 及び s_2 のリンクの重みは、以下のようになる。

【0068】

$$lw(p_1, q) = 1 / \{1 \times (1 + 1 + 1)\} = 1 / 3$$

$$lw(p_2, q) = 1 / \{2 \times (1 / 2)\} = 1$$

$$lw(p_3, q) = 1$$

$$lw(s_1, p_3) = lw(s_2, p_3) = 1 / 3$$

故に、(1) 式及び上述の計算からリンク先の多いページ p_1 のリンクの重み $lw(p_1, q)$ は小さくなることが分かる。また、同様に(1) 式及び上述の計算からURL類似度が大きいほどリンクの重みは小さくなることも分かる。

【0069】

また、ページ q のリンク重要度 W_q は以下のようになる。

$$\begin{aligned} W_q &= C_q + \{lw(p_1, q) \times W_{p_1} + lw(p_2, q) \times W_{p_2} + lw(p_3, q) \times W_{p_3}\} \\ &= C_q + \{(W_{p_1} / 3) + W_{p_2} + W_{p_3}\} \end{aligned}$$

$$W_{p_1} = C_{p_1}$$

$$W_{p_2} = C_{p_2}$$

$$W_{p3} = C_{p3} + \{lw(s1, p3) \times W_{s1} + lw(s2, p3) \times W_{s2}\}$$

$$= C_{p3} + (W_{s1} + W_{s2}) / 3$$

故に、ページ p_1 及び p_2 と比べてより多くのページからリンクされているページ p_3 のリンク重要度 W_{p3} は大きくなっている。また、ページ q のリンク重要度 W_q をみれば、重要なページと思われるから、URL 類似度の低いページ、すなわちリンク重み $lw(p3, q) = 1$ 、いいかえるとリンク重み（リンク重みの最大値は $=$ ）が大きいページほど、リンク重要度が大きくなることが分かる。また、（2）式及び上述の計算から、URL の類似している同じサイト内のページからのリンクの重みを、他の URL の類似していないページからのリンクの重みよりも軽く扱うことが分かる。これにより、大量のページを抱えるサイトのページが検索結果の中に多く出てくるという問題を解決できていることが分かる。

【0070】

図8は、（1）式及び（2）式の示す内容を模式的に示す。図8においても、図7と同様に、まるは各ページ、矢印はリンク関係、矢印の方向はリンクの方向、矢印の太さがリンクの重みを示す。さらに、図8において網掛されているまるはURL 類似度が高いページを示す。図8（a）及び8（b）において、ページ q はページ p_1 、 p_2 及び p_3 からリンクされている。さらに、図8（b）において、ページ q はページ p_1 、 p_2 及び p_3 のURL は類似しており、URL 類似度 $sim(pi, q) = n + 1$ である（ n は整数）。図8（a）及び図8（b）それぞれについて（1）式及び（2）式を適用する。図8（a）の場合は以下のようになる。

【0071】

各ページのリンクの重み

$$lw(pi, q) = 1 / sim(pi, q) = 1 \quad (\text{URL が非類似})$$

ページ q のリンク重要度 W_q

$$W_q = C_q + (W_{p1} + W_{p2} + W_{p3})$$

図8（b）の場合は以下のようになる。

【0072】

各ページのリンクの重み

$$lw(pi, q) = 1 / sim(pi, q) = 1 / (n + 1) \quad (URL \text{ が類似})$$

ページ q のリンク重要度 W_q

$$W_q = C_q + (W_{p1} + W_{p2} + W_{p3}) / (n + 1)$$

故に、図 8 (a) 及び 8 (b) それぞれの計算結果を比較すると、URL 類似度 $sim(p, q)$ が高い場合には、被リンク数が多くてもページ q のリンク重要度 W_q が小さくなることが分かる。延いては、URL 類似度を導入することにより、大量のページを抱えるサーバ (サイト) 等が単にページの量が多いというだけで重要度が高いことになるという問題を解決していることが分かる。

【 0 0 7 3 】

次に、(1) 式及び (2) 式中のページ p と q の URL 類似度 $sim(p, q)$ について説明する。URL 類似度は、リンク重要度付与器 2 1 に備えられた URL 類似度計算器 2 7 により算出される。

【 0 0 7 4 】

一般に、ページの URL は、サーバアドレス、パス、ファイル名の三種類の情報から構成される。例えば、WWW ページの URL、

<http://www.flab.fujitsu.co.jp/hypertext/news/1999/product1.html> は、サーバアドレス (www.flab.fujitsu.co.jp)、パス ([hypertext/news/1999](http://www.flab.fujitsu.co.jp/hypertext/news/1999))、ファイル名 ([product1.html](http://www.flab.fujitsu.co.jp/hypertext/news/1999/product1.html)) の 3 種類の情報から構成される。

【 0 0 7 5 】

また、サーバアドレスは、さらに “.” により階層化されており、後ろに行くにしたがって、段々広くなる。例えば、サーバアドレスが www.flab.fujitsu.co.jp であれば、後ろから、日本 (jp)、会社 (co)、富士通 (fujitsu)、研究所 (flab)、マシン (www) という階層を表している。

【 0 0 7 6 】

本実施形態では、与えられた 2 つのページ p 及び q の URL 類似度を、上記の三種類の組合せにより定義する。類似度 $sim(p, q)$ として、例えば、以下に述べるドメイン類似度 $sim_domain(p, q)$ 及び融合類似度 $sim_merge(p, q)$ が考えられる。

【 0 0 7 7 】

ドメイン類似度 $\text{sim_domain}(p,q)$ は、ドメインの類似に基づいて算出される。ドメインとは、サーバアドレスの後半部分であり、会社や組織を表す。サーバアドレスが .com、.edu、.org 等で終わる米国サーバの場合はサーバアドレスの後ろから2つめまで、サーバアドレスが .jp、.fr 等で終わる他国のサーバの場合はサーバアドレスの後ろから3つめまでがドメインに相当する。例えば、www.fujitsu.com のドメインは fujitsu.com であり、www.flab.fujitsu.co.jp のドメインは fujitsu.co.jp である。

【0078】

ページ p とページ q のドメイン類似度は以下の(3)式により定義される。

$$\begin{aligned}\text{sim_domain}(p,q) &= 1/\alpha && (p, q \text{ が同一ドメインの場合}) \\ &= 1 && (p, q \text{ が異なるドメインの場合})\end{aligned}$$

ここで、 α は定数で、0より大きく1より小さい実数値を取るとする。図9はインターネットから収集した約300万URL間のリンク関係にドメイン類似度 $\text{sim_domain}(p,q)$ の概念を導入してリンク重要度を計算した場合を示す。図9において、横軸が、リンク重要度が高い順から上位何ページかを示す数、縦軸が上位ページの中に含まれる異なるドメインを持つページの数を示し、系列1から5は、それぞれ順に $\alpha = 0.1, 0.2, 0.3, 0.5, 0.7$ 及び 1.0 の場合を示す。図9に示すように、リンク重要度が高い上位10万件のページ中に含まれる異なるドメインを持つページ数は、 $\alpha = 1$ の場合 (URL類似度を導入しない従来の場合に相当する) は4000件であり、 $\alpha = 0.1$ の場合は5500件である。従って、 α が小さくなるほど、異なるドメインを持つページがリンク重要度の上位に上がってくることが分かる。これは、 α が小さいほど類似度 $\text{sim_domain}(p,q)$ が大きくなり、URL類似度 $\text{sim_domain}(p,q)$ が大きいほど、リンクの重み $lw(p,q)$ が小さくなり、延いては、リンク重要度 W_q が小さくなるため、URL類似度が大きくなるほどページには小さいリンク重要度が付与されるからである。つまり、 $\text{sim_domain}(p,q)$ の概念を導入することにより、異なるドメインを持つページが検索されやすくなっている、言い換えると、同じドメインを持つページは検索されにくくなっていることが分かる。

【0079】

$\text{sim}(p,q)$ として、前述の三種類の情報を融合した類似度 $\text{sim_merge}(p,q)$ を次のように定義する。

$\text{sim_merge}(p,q) = (\text{サーバアドレスの類似度}) + (\text{パスの類似度}) + (\text{ファイル名の類似度})$

以下、右辺の各項の算出方法について説明する。

【 0 0 8 0 】

サーバアドレスの類似度は、アドレスの階層を後ろから見ていき、 n レベルまで一致した場合、類似度を $1 + n$ とする。例えば、`www.fujitsu.co.jp` と `www.flab.fujitsu.co.jp` は 3 レベルまで一致しているので 4 となる。`www.fujitsu.co.jp` と `www.fujitsu.com` は 1 レベルも一致していないので（一致 0 レベル）、類似度は 1 である。

【 0 0 8 1 】

パスの類似度は、先頭からパスの“/” で区切られた要素毎に比較し、一致したレベルまでを類似度とする。例えば、`/doc/patent/index.html` と `/doc/patent/1999/2/file.html` とは、2 レベルまで一致しているので類似度は 2 である。

【 0 0 8 2 】

ファイル名の類似度は、ファイル名が一致する場合、類似度 1 とする。

上記は、以下のような考え方に基づいている。

1. 往々にして、同じような文書を同一ディレクトリに入れるため、同一サーバでパスも同じ URL は内容が似ていることが多い。
2. アクセスを分散させるために設けられるミラーサイトは類似度が高い。これらは、サーバアドレス部分だけが異なり、残りのパスやファイル名は同じ場合が多い。
3. サーバアドレス、パス、ファイル名が全てことなる URL は、類似度が低い。

【 0 0 8 3 】

この $\text{sim_merge}(p,q)$ によっても、URL が似通ったページが検索されることを防ぐことができる。従って、 $\text{lw}(p,q)$ の中に $\text{sim}(p,q)$ 又は $\text{diff}(p,q)$ という概念を導入することにより、大量のページを抱えるサーバ（サイト）や個人が単に

量が多いというだけでリンク重要度が高いことになるという従来の問題を解消することができる。

【 0 0 8 4 】

なお、上述のリンク重要度 W_p は、後述の関連度算出にも用いられる。

以下、文書検索装置のキーワード文書関連度計算器 2 3 が関連度を算出する手順について説明する。

【 0 0 8 5 】

まず、キーワードから文書への索引を作るにあたっては、キーワードと文書への関連度が必要である。関連度については、以下のように考える。

1. キーワードを多く含んでいる、つまりキーワードを含む度合いが高い文書ほど関連度が大きい。
2. 重要度の高い文書ほど関連度が大きい。
3. あるキーワードの関連文書は一定個数以下であることが望ましい。（一つのキーワードで、1,000 個もの関連文書が得られるのはまずい）。

【 0 0 8 6 】

本実施形態ではあるキーワードの関連文書は一定個数以内にするために、上記の考え方に加え、以下の考え方を導入する。

4. 利用者のアクセスログ解析による関連度：一定期間内に当該キーワードからよくアクセスされた文書ほど関連度が高い。
5. 文書のリンク重要度による関連度：当該キーワードを含む文書のうち、リンク重要度の高いものは関連度が高い。

【 0 0 8 7 】

上述の考え方を導入して、あるキーワード w に対してページ p の関連度を以下の (3) 式で与える。

$$Rel(p, w) = TF(p, w) * \log W_p * \log(AC(p, w) + 2) \quad \cdots (3)$$

ここで、 $TF(p, w)$ は p における w の出現個数、

W_p は p のリンク重要度であり、上述の (2) 式の W_p に相当する。

【 0 0 8 8 】

$AC(p, w)$ は、一定期間（例えば 1 月間又は 1 週間）の間にキーワード w から p にアクセスされた回数である。

各キーワードに対して、上記 $Rel(p, w)$ の値の大きいものから一定個数をそのキーワードの関連ページとする。

【 0 0 8 9 】

キーワードを含む度合いに加えて、URL類似度という概念を導入したリンク重要度 $W(p)$ 及び利用者のアクセスログ解析を関連度の算出に用いている。これにより、ページの関連度を高くするための条件が多くなるため、意図的に第三者があるページの関連度を高くするようにページ内容进行操作することはさらに困難になる。

【 0 0 9 0 】

以下、索引生成器 2 4 が生成する索引、つまり、検索のためのキーワードを選択するインターフェースについて説明する。本実施形態の選択インターフェースは、キーワードの読みの断片を順次クリックするだけでキーワードに辿りつけるインタフェースであり、以下の特徴がある。

- ・一画面には、キーワードの読みの断片（文字または文字列）と、それまで選択したキーワード読みに相当するキーワードの一部が表示される。
- ・利用者は、画面上のキーワードの読み断片（文字または文字列）を次々クリックすることで、キーワードに辿りつける。
- ・1画面に表示されるキーワードは、一定個数以内に制限することができる。

【 0 0 9 1 】

前述の従来技術では、常に1「文字」をクリックして読みを選択しているが、本実施形態では必ずしも文字でなく、「文字列」の場合もあり得る。これにより、キーワードに辿りつくまでの無駄なクリックを減らすことが可能である。さらに、1画面に表示されるキーワードは一定個数以内であるため、キーワードを選択することが用意となる。また、1画面に表示されるキーワードは一定個数以内であることは、表示装置の画面の大きさが制限されている携帯電話のような移動端末でのキーワード選択に有用である。これらを実現するために、索引生成器 2 4 において以下のようにしている。

1. キーワードの読み（またはスペル）の統一。場合によっては「ー」（長音）の除去や「い」のような拗撥音を「あ」「い」のような表記に統一する。

2. キーワードとその読み（またはスペル）から、読み文字をノード、キーワード集合をリーフとする有効グラフ（文字列グラフ）を作る。

3. 上記グラフで、以下の縮退操作を行う（下を参照）。

【0092】

(a) リーフへのパスの縮退

(b) 中間パスの除去

(c) 子ノードのキーワードを親ノードにくりこみ、子ノードの削除

4. 縮退後のグラフを元に、読み入力インタフェースを作成する。

【0093】

以下、索引生成器24においてキーワード文字列グラフ作成以下の手順をどのように行っているか説明する。

キーワード文字列グラフとは、キーワードの読みを表現する有向ラベル付きグラフである。図10は、キーワード文字列グラフの一例を示す。

【0094】

キーワード文字列グラフは、

$(N, C, KW, t, nk, yomi)$

の6つ組で表現される。ここで、

N はノードの集合、

C はかな文字の集合、

KW はキーワードの集合、

t は $N * C^+ \rightarrow N$ ノード遷移関数。また、 C^+ はラベル、つまり、1以上のかな文字の並び（図9（a）等の文字列グラフでは実線矢印で示される）

nk は $N \rightarrow W^+$ ノードに割り当てられたキーワード（図9（a）等では点線で示される）

$yomi$ は $N \rightarrow C^+$ ノードの「読み」である。

【0095】

図10（a）を例にとると以下のようなになる。なお、 $yomi$ は自明なので省略）

$N = \{\text{トップ}, \text{「あ」}, \text{「あい」}, \text{「あいぼ」}, \text{「あいぼう」}, \text{「あいぼり」}, \text{「あお」}, \text{「あおぞ」}, \text{「あおぞら」}\}$

$C = \{\text{' あ' }, \dots, \text{' ん' }\}$

$KW = \{\text{' 青' }, \text{' 蒼' }, \text{' 青空' }, \text{' アイボリー' }\}$

$t(\text{トップ}, \text{あ}) = \text{「あ」},$

$t(\text{あ}, \text{い}) = \text{「あい」},$

$t(\text{あ}, \text{お}) = \text{「あお」},$

$t(\text{あい}, \text{ぼ}) = \text{「あいぼ」},$

$t(\text{あいぼ}, \text{う}) = \text{「あいぼう」},$

$t(\text{あいぼ}, \text{り}) = \text{「あいぼり」},$

$t(\text{あお}, \text{ぞ}) = \text{「あおぞ」},$

$t(\text{あおぞ}, \text{ら}) = \text{「あおぞら」},$

$nk(\text{あいぼう}) = \{\text{' 相棒' }\}$

$nk(\text{あいぼり}) = \{\text{' アイボリー' }\}$

$nk(\text{あお}) = \{\text{' 青' }, \text{' 蒼' }\}$

$nk(\text{あおぞら}) = \{\text{' 青空' }\}$

索引生成器24はキーワードと読みが与えられた場合、それらに基づいて初期キーワード文字列グラフを生成する。図11は初期キーワード文字列グラフの生成手順を示す。図11を用いて索引生成器24における初期キーワード文字列グラフの生成手順を説明する。なお、図12は、初期キーワード文字列グラフの生成手順を実現するアルゴリズムの一例を示す。

【0096】

まず、キーワードの集合KWを作成する(ステップS11)。次に、作成された集合KW内が空か否か判定する(ステップS12)。集合KW内が空である場合(ステップS12、yes)、文字列を作成する必要はないため、処理を終了する。集合KW内が空でない場合(ステップS12、no)、次のステップに進む。

【0097】

次に、集合KWからあるキーワードuを抜き出す(ステップS13)。そして

、このキーワード u の読み $yomi(u)$ 及び、この読み $yomi(u)$ のノードである $nk\{yomi(n)\}$ を設定し、このノード $nk\{yomi(n)\}$ を末端ノードとして追加する（ステップ S 1 4）。

【0098】

ステップ S 1 4 の処理をキーワード u の文字列長分繰り返したか否か、つまり、キーワード u が空になったか否か判定する（ステップ S 1 5）。キーワード u が空になったか場合（ステップ S 1 5、yes）、キーワード u については必要な処理は終了したため、ステップ S 1 2 へ戻り集合 KW から別のキーワードをとってステップ S 1 3 以下の手順を繰り返す。キーワード u が空になっていない場合（ステップ S 1 5、no）、キーワード u の末尾文字を切り取る（ステップ S 1 6）。続いて、ノードを 1 つ前の親ノードに遷移して追加する（ステップ S 1 7）。さらにキーワード u の切り取られた末尾文字の前の文字に注目し（ステップ S 1 8）、ステップ S 1 5 に戻る。

【0099】

この手順により、集合 nk としてノードに割り当てられたキーワードリストが得られ、 t としてあるノードの下位に位置する下位ノードのリストが得られる。

図 1 6（a）は上述の手続により作成された初期キーワード文字列グラフを示す。図 1 6（a）は、

蒼：あお、

青：あお、

青空：あおぞら、

相棒：あいぼう、

アイボリー：あいぼり

というキーワードと読みデータから生成した初期キーワード文字列グラフである。また、同様に、図 9（a）は図 1 1 に示す `init_kw_graph()` アルゴリズムにおいて

@KW = {蒼、青、青空、相棒、アイボリー}

$yomi\{\text{蒼}\} = \text{あお}$ 、 $yomi\{\text{青}\} = \text{あお}$ 、 $yomi\{\text{青空}\} = \text{あおぞら}$ 、 $yomi\{\text{相棒}\} = \text{あいぼう}$ 、 $yomi\{\text{アイボリー}\} = \text{あいぼり}$ 、

としたものである。

【0100】

索引生成器24は、初期キーワード文字列グラフを作成すると、次に、文字列の縮退を行う。以下、文字列の縮退について説明する。縮退は、

1. 中間ノードの縮退
2. 末端ノードのくりこみ

の2種類の操作からなる。

【0101】

まず、索引生成器24における中間ノードの縮退処理について説明する。図13は中間ノードの縮退処理の手順を示す。以下、図13を用いて中間ノードの縮退処理の手順について説明する。なお、図14は、中間ノードの縮退処理を実現するアルゴリズムの一例を示す。

【0102】

まず、ノードの集合Nを作成する（ステップS21）。続いて、集合N内が空か否か判定する（ステップS22）。集合N内が空である場合（ステップS22、yes）、ノードを縮退させる必要はないため、処理を終了する。集合N内が空でない場合（ステップS22、no）、集合N内のノードnを取得する（ステップS23）。取得したノードnについて、そのノードnに続く後継ノードが1つしかなく、かつノードn内にキーワードもないという2つの条件を満たすか否か判定する（ステップS24）。2つの条件を満たす場合（ステップS24、yes）、ノードnは縮退可能であるため、ノードnをキーワード文字列グラフから削除し、ステップS22に戻る（ステップS25）。2つの条件の両方を満たさない場合（ステップS24、no）、ノードnは縮退できないので、ノードnを削除せずにステップS22に戻る。

【0103】

上述のようにしてキーワード文字列グラフにおいて、「ノードに割り当てられたキーワードがない」又は、「後継ノード（子ノード）が1つのみ」であるという2つの条件を満たす中間ノードは縮退させることとしている。図10（a）の初期キーワード文字列グラフでは、ノード「あい」及びノード「あおぞ」が「ノ

ードに割り当てられたキーワードがない」又は、「後継ノード（子ノード）が1つのみ」であるという2つの条件を満たす中間ノードに相当する。図10（b）は、図10（a）に示す初期キーワード文字列グラフの中間ノードを縮退させた結果を示す。図10（b）において、中間ノード「あい」及び「あおぞ」は削除されている。また、同様に、図14に示すproc_shrink_middle()アルゴリズムにおいて、

```
t {''} = あ +
t {あ} = あいぼ + あお +
t {あいぼ} = あいぼう + あいぼり +
t {あお} = あおぞ +
nk {あいぼう} = 相棒 +
nk {あいぼり} = アイボリ +
nk {あお} = 青 + 蒼 +
nk {あおぞら} = 青空 +
```

となる。

【0104】

次に、索引生成器24における末端ノードの縮退処理について説明する。図15は末端ノードの縮退処理の手順を示す。以下、図15を用いて末端ノードの縮退処理の手順について説明する。なお、図16は、末端ノードの縮退処理を実現するアルゴリズムの一例を示す。

【0105】

まず、全てのノードの集合Nを作成する（ステップS31）。次に、ノードをそのノードに属するキーワード数の順にソートする（ステップS32）。整数iを1に設定する（ステップS33）。次に、整数iが集合N内のノードの数よりも少ないか否か判定する（ステップS34）。整数iが集合N内のノードの数よりも少ないのではない場合（ステップS34、no）、さらに、末端ノードの縮退があるか否か判定する（ステップS35）。末端ノードの縮退がない場合（ステップS35、no）、処理を終了する。末端ノードの縮退がある場合（ステップS35、yes）、ステップS33に戻る。

【0106】

整数 i が集合 N 内のノードの数よりも少ない場合（ステップ S 3 4、*yes*）、集合 N の i 番目のノード n を取得する（ステップ S 3 6）。取得されたノード n が末端ノードであるか否か判定する（ステップ S 3 7）。取得されたノード n が末端ノードである場合（ステップ S 3 7、*yes*）、ステップ S 3 8 に進む。取得されたノード n が末端ノードでない場合（ステップ S 3 7、*no*）、ノード n は縮退すべきノードではないため、整数 i を 1 つインクリメントし（ステップ S 4 1）、ステップ S 3 4 に戻る。

【0107】

次に、ノード n の親ノード p を取得する（ステップ S 3 8）。次に、親ノード p に属するキーワードの数と子ノード n に属するキーワードの数の和が予め与えられたキーワードの数の上限を超えるか否か判定する（ステップ S 3 9）。

【0108】

親ノード p に属するキーワードの数と子ノード n に属するキーワードの数の和が上限値を超えない場合（ステップ S 3 9、*yes*）、子ノード n を削除し（縮退させる）、ノード n に属するキーワードを親ノード p に属するキーワードとする（ステップ S 4 0）。続いて、整数 i を 1 つインクリメントして（ステップ S 4 1）、ステップ S 3 4 へ戻る。

【0109】

親ノード p に属するキーワードの数と子ノード n に属するキーワードの数の和が上限値を超える場合、子ノード n を縮退させると親ノードのキーワード数が多すぎることになるため、子ノード n を縮退させずに、ステップ S 4 1 に進む。

【0110】

このようにして末端ノードにあるキーワード情報を、一つ上の親ノードに移すことにより、木（ノードの連鎖）の深さを減らし、利用者が少ないクリックでキーワードに到達できるようになる。しかし、あまり沢山親ノードに子ノードのキーワードを移しすぎると、一つのノードに大量のキーワードが割り当てられることになるので、利用者がその中からキーワードを選択することが面倒になってしまう。この問題を回避するため、パラメータ `words __max` を与え、1 ノードに属

するキーワードがそのパラメータ `words __max` より少なくなるようにする。

【0111】

図10(c)は、図10(b)のキーワード文字列グラフをパラメータ `words __max = 4` の場合に末端ノード縮退処理した結果を示す。図10(b)において、末端ノード「あいぼう」及び「あいぼり」はそれぞれ1つのキーワードしか持たない。末端ノード「あいぼう」及び「あいぼり」の親ノード「あいぼ」は2つの子ノードを持ち、キーワードはない。従って、親ノード「あいぼ」と子ノード「あいぼう」及び「あいぼり」のキーワードの和は `words __max = 4` より少ない。故に、子ノード「あいぼう」及び「あいぼり」は縮退可能であるため、図10(c)において子ノード「あいぼう」及び「あいぼり」は削除され、これら子ノードのキーワードは親ノード「あいぼ」に移されている。また、同様に、図16に示す `proc__shrink__leaf()` アルゴリズムにおいて、

t {'} = あ +

t {あ} = あいぼ + あお +

t {あお} = あおぞ +

n k {あいぼ} = 相棒 + アイボリー +

n k {あお} = 青 + 蒼 +

n k {あおぞら} = 青空 +

となる。

【0112】

図17は、末端ノード縮退処理をしたキーワード文字列グラフの更なる一例を示す。図17において、親ノード「かいせ」には3つの末端ノード「かいせい」、「かいせき」及び「かいせつ」がある。末端ノード「かいせい」及び「かいせき」に属するキーワードはそれぞれ「快晴」及び「解析」であるため、末端ノード「かいせい」及び「かいせき」は縮退可能である。また、末端ノード「かいせつ」に属するキーワードは「解説」及び「開設」の2つであるため、末端ノード「かいせつ」も縮退可能である。従って、図17の場合は2通りの縮退の方法が考えられる。しかし、前者の方が総ノード数を減らすことができる。本実施形態において、末端ノードをそのノードに属するキーワード数の順にソートしている

。これにより、前者のような効率的な縮退をすることを可能にしている。

【0113】

次に、索引生成器が作成する索引の例について図18から図25を用いて説明する。図18は索引トップ画面から索引中間画面及びキーワード情報画面を経て分所ページに至るまでの索引画面の遷移を示す。図18を用いて表示装置に表示される索引画面の遷移について説明する。図18において、まず、索引トップ画面が表示される。利用者が索引トップ画面からキーワードの読み又はスペルの最初の部分を選択すると、索引中間画面に遷移する。利用者が索引中間画面からキーワードの読み又はスペルの次の部分を選択すると、さらなる索引中間画面に遷移する。この選択を繰り返した結果、利用者が検索を希望するキーワードに辿りつき、そのキーワードを選択すると、キーワード情報画面に遷移する。利用者がその他のキーワードを選択すると、さらなるキーワード情報画面に遷移する。利用者が閲覧を希望する文書ページのタイトルを選択すると、その文書ページへリンクし、その文書ページに遷移する。なお、利用者の選択操作は、マウスによるクリックやペン等のポインティングにより行うようにすることも可能である。各画面は、例えばハイパーテキストで作成されてもよい。

【0114】

図19は索引トップ画面の一例を示す。索引トップ画面には、図4に示す索引情報テーブル61のトップ以下の文字（列）が表示される。図19において、50音、アルファベット及び0から9までの数字が表示されている。利用者は検索したいキーワードの読みまたはスペルの最初をクリックすることで先の画面に進むことができる。

【0115】

図20は索引トップ画面のさらなる例を示す。図20において、キーワードの読み付けをする際に表記統一した又は・及びノードの縮退をした結果、例えば、50音の「だ」行の「ぢ」及び「づ」が索引から除かれていることが分かる。同様に、アルファベットの「Y」及び「Z」が除かれていること等が分かる。

【0116】

図21は索引中間画面の一例を示す。図21は、索引トップ画面において「あ

」を選択した場合の索引中間画面である。上半分は「あ」の後に継続する文字列、下半分にはその他のキーワードが表示されている。索引中間画面は、図4の索引情報テーブル61及び図3のキーワード対応テーブル51（キーワードから、キーワードIDを得る）から生成できる。

【0117】

図21において、「あ」に継続する文字列として「いぼ」、「え」「お」などがある。例えば、「いぼ」を選択すると「あいぼ」に移る。その他のキーワードとして、「愛」「愛犬」などが一定個数（例えば20以内）表示されている「あ」から始まっていて、継続する読みが選べないキーワードは全てここに表示されている。逆にキーワードも表示されず、継続する読みもない場合には、そのキーワードは索引に含まれていないことがわかる。

【0118】

図22は、索引中間画面の更なる例を示す。図22は、索引トップ画面において「い」を選択した場合の索引中間画面である。上半分は「い」の後に継続する文字列、下半分にはその他のキーワードが表示されている。図22において、「あ」に継続する文字列として「いぼ」、「えろ」などがある。その他のキーワードとして、「イオン」「イネーブル」などが一定個数（例えば20以内）表示されている。

【0119】

図23は、索引中間画面のさらなる例を示す。図23において、「いべんと」が選択された場合の索引中間画面が示されている。ノード「イベント」の子ノードはないため、「いべんと」の後に継続する文字列は表示されず、キーワードが表示されている。利用者は、この画面から入力したいキーワードを選択する。この画面に含まれていない「いべんと」を含むキーワードは索引に含まれていないことが分かる。

【0120】

従来技術において読みは1文字ずつしか選択できないため、互い文字列を入力するためには何度も選択操作を繰り返す必要があった。しかし、本実施形態によれば、ノードを縮退させるため、継続する文字列を必ずしも1文字ずつ選択する



必要はない。例えば図20の文字列「いぼ」又は図21の文字列「えろ」のように、2文字を一度に選択することも可能である。従って、利用者が選択する回数を低減することができる。

【0121】

また、さらに、継続する読みが選べないキーワードは全てここに表示されている。逆にキーワードも表示されず、継続する読みもない場合には、そのキーワードは索引に含まれていないため、利用者が1文字ずつ選択していった最後の1文字をするときになって初めて入力しようとしていたキーワードが索引に含まれていないことが分かるといった問題は生じないことになる。

【0122】

また、さらに、末端ノードの縮退の際に、パラメータword#maxとして設定された値以上のキーワードは表示されないため、利用者はキーワードの選択の際に表示されたキーワードから選択したいキーワードを探し出すことが比較的楽におこなうことができる。また、表示されるキーワードが一定個数以内にかぎられていることは、表示画面の大きさが限られた携帯電話のような移動端末におけるキーワード選択において有用である。

【0123】

さらにまた、検索のインタフェースとして、決まったキーワードをセットに、なるべく少ない労力で辿り着かせることを可能にし、従来技術におけるかな漢字変換とは次のような点が異なる。

- ・変換キー操作は不要
- ・変換したいキーワードの読みの全てを入れなくても、そのキーワードを特定できる最低限の情報さえ与えられれば良い。

【0124】

従って、例えば、キーワードセットに、「なれ」から始まる語として「ナレッジマネジメント」しか入っていない場合、利用者が「な」「れ」と指定したら、それだけで「ナレッジマネジメント」を表示する。

【0125】

図24は、キーワード情報画面の一例を示す。中間画面でキーワード「圧縮」



をクリックした時のキーワード情報画面が示されている。図 2 4 において、まず、代表語と同義語（「圧縮」「Compress」）とが表示される。これは、図 3 のキーワードテーブル 5 1 とキーワード対応テーブル 5 2 の情報から得ることができる。画面右上に「あ」とあるのは、辿ってきたパスを示す。これにより、前に戻ったり、ハイパーテキストで迷子になるのを防ぐことができる。また、画面には各文書のタイトル、リンク情報、他のキーワードが表示されている。優先度の高い順に一定個数（例えば 20 文書）以内が表示されるので、利用者はこの中から文書を選択するのが困難ということはない。文書 ID のリストは、図 4 の関連文書テーブル 6 2 から得ることができる。各文書 ID についての情報は、図 2 の文書情報テーブル 4 1 に入っている。他のキーワードは、図 4 の関連キーワードテーブル 6 3 から得る。利用者は閲覧を希望する文書情報を選択すると、その文書へリンクして文書を画面に表示させることができる。

【 0 1 2 6 】

図 2 5 は、キーワード情報画面のさらなる例を示す。中間画面でキーワード「イベントカレンダー」をクリックした時のキーワード情報画面が示されている。画面右上に「トップ」－「イ」－「イベント」とあるのは、辿ってきたパスを示す。それぞれのパスを選択することにより、利用者は先に見た画面に戻ることができる。

【 0 1 2 7 】

図 2 6 は第 2 実施形態に係わるイントラネット文書検索装置の構成を示す。図 2 6 において、図 1 に示す第 1 実施形態の構成に加えて、収集装置 8 1 及び同義語辞書 8 2 を更に備える。収集装置 8 1 はイントラネット（又はインターネット）から大量の文書を収集する、例えば Web ロボットである。同義語辞書（同義語データ）8 2 は同キーワード対応テーブルの情報の一部を格納する。なお、入力装置及び出力装置は、例えば、WWW ブラウザ 8 3 であってもよい。

【 0 1 2 8 】

収集装置 8 1 はネットワークから文書を自動で収集し、処理装置 1 1 のキーワード抽出器 2 2 は同義語辞書 8 2 を用いて、収集されたページからキーワードを抽出し、ページ内のキーワード出現頻度を集計する。つまり、同義語辞書 8 2 を

用いて重要度の高い文書を自動選別する。これにより、イントラネット（又はインターネット）の大量の文書を自動選別することを可能にする。

【 0 1 2 9 】

図 2 7 は第 3 実施形態に係わる特定タイプ文書のイントラネット文書検索装置の構成を示す。図 2 7 において、図 2 6 に示す第 2 実施形態の構成に加えて、処理装置 1 1 内に文書タイプ判別器 9 1 を備える。文書タイプ判別器 9 1 はイントラネット（及び・又はインターネット）から収集された文書データのリンク関係及び URL に基づいて、その文書データの文書タイプを判別する。より具体的には、文書タイプ判別器 9 1 はリンク重要度付与器 2 1 内で算出された URL 類似度とリンク重要度付与器 2 1 により抽出されたリンク関係が示すリンク／被リンク数に基づいて、（その内容を理解せずに）文書データのコンテンツのタイプを判別する。文書タイプ判別器 9 1 における文書タイプの判別は以下のように行う。

1. URL 類似度が一定以上の文書データへのリンクを一定数以上持つページはリンク集である。
2. URL 類似度が一定以下の文書データへのリンクを一定数以上持つページはメニュー（エントリ）ページである。
3. URL 類似度が一定以下の文書データから一定数以上参照されているページは、メニュー（エントリ）ページである。
4. それ以外で、URL 類似度が一定以上文書データへのリンクが一定数以下のページはコンテンツページである。

【 0 1 3 0 】

このように判別することにより、文書タイプ判別器 9 1 は、かなりの確率で文書データ（WWW ページ）の文書タイプ（例えば、メニューページ、リンク集、コンテンツページ）を分けることができる。

【 0 1 3 1 】

文書タイプ判別器 9 1 は、文書タイプを判別し、判別された文書タイプ 9 2 をキーワード文書関連度計算器 2 3 に出力する。キーワード文書関連度計算器 2 3 は、判別された文書タイプ 9 2 に基づいて特定タイプの文書データを選択し、選

択された文書データについてリンク重要度、ページキーワード及びアクセスログに基づいて文書関連度を計算する。例えば、キーワード文書関連度計算器 2 3 は文書タイプがコンテンツページであると判別された文書データを選択して、これらのコンテンツページについての関連度を計算するようにしてもよい。

【 0 1 3 2 】

このように、図 2 6 のイントラネット文書検索装置は、索引に掲載するページをここで得た文書タイプ判別器の判別したタイプのいずれかに限定することで、より質の高い文書整理を行うことが可能となる。

【 0 1 3 3 】

図 2 8 は第 4 実施形態に係わるリンク集生成システムの構成を示す。図 2 8 において、リンク集生成システムは、収集装置 1 0 1、処理装置 1 0 2 及び入出力装置 1 0 7 を備える。収集装置 1 0 1 はインターネット（又は・及びイントラネット）から収集装置 8 1 は大量の文書データを収集する、例えば Web ロボットである。処理装置 1 0 2 は、リンク重要度付与器 2 1、URL 文字列判別器 1 0 3、索引生成器 2 4 及び WWW サーバ 1 0 6 を備える。リンク重要度付与器 2 1 は URL の類似度とリンク関係を元に、文書のリンク重要度を算出し、算出されたリンク重要度 3 1 を索引生成器 2 4 へ出力する。

【 0 1 3 4 】

URL 文字列判別器 1 0 3 は URL の文字列上の特徴に基づいて、収集された文書データのコンテンツを（その内容を理解せずに）判別する。URL 文字列判別器 1 0 3 は、URL の文字列上の特徴から例えば、以下のように文書データのコンテンツを判別する。

1. URL の文字列に Y2K, y2k, year2000 を含むような文書データは 2000 年問題の関連ページである。
2. URL の文字列に news, release, press を含み、その後で数字列（往々にして日時の情報）を持つ文書データはニュース（プレス）リリースのページである。
3. URL の文字列に java, JAVA を含むページは Java 関連ページである。
4. URL の文字列に download, dwnload, dwnld を含むページは、ソフトウェア

ダウンロード関係のページである。

5. URLの文字列にLINUX, linux, Linux を含むページは、Linux 関連ページである。

【0135】

このように判別することにより、URL文字列判別器103は特定のURLを判別し、判別された特定URL集合104を索引生成器24に出力する。索引生成器24は、リンク重要度31に基づいて、特定URL集合104をリンク重要度の大きい順に並べて、上位から一定個数を取り出し、上位URLからリンク集を作成し、リンク集105としてWWWサーバ106に出力する。なお、文字列判別器103で得られたURLが少数の場合には、リンク関係を調べてそれらのリンク元がよく参照（リンク）する他のページを加えて数を増やしてもよい。これは、しばしばWWWで似たようなページは、似たようなリンク集から同時に参照されるためである。WWWサーバ106はリンク集を利用者に提供し、利用者はWWWブラウザ107を介してリンク集を表示させ、入力指示を行う。

【0136】

このようにして、URL文字列に基づいて、文書ページのコンテンツをみないでコンテンツを判別し、この判別結果を用いてリンク集を作成する。これにより、コンテンツの内容にそった質の高いリンク集を簡便に作成することが可能となる。

【0137】

図29は第5実施形態に係わるリンク集生成システムの構成を示す。図29のリンク集生成システムは、図28のリンク集生成システムに文書タイプ判別器111を更に備える構成となっている。文書タイプ判別器111の機能及び動作は、図27を用いて説明した第3実施形態に係わる文書検索装置の文書タイプ判別器91と同様である。

【0138】

収集装置101はインターネット（又は・及びイントラネット）から収集装置81は大量の文書データを収集し、リンク重要度付与器21はURLの類似度とリンク関係を元に文書のリンク重要度を算出し、リンク重要度31を索引生成期

24へ出力する。URL文字列判別器103はURLの文字列上の特徴に基づいて、特定のURLを判別し、判別された特定URL集合104を索引生成器24へ出力する。文書タイプ判別器111は、URL類似度とリンク／被リンク数に基づいて、（コンテンツを見ずに）文書データの文書タイプを判別し、判別された文書タイプ112を索引生成器24に出力する。

【0139】

索引生成器24は、文書タイプ112に基づいて特定URL集合104の中から特定の文書タイプの文書データを選び出す。続いて、選ばれた文書データのリンク重要度31に基づいて、選ばれた文書データをリンク重要度順の大きい順に並べて、上位から一定個数を取り出し、上位URLからリンク集を作成し、リンク集105としてWWWサーバ106に出力する。WWWサーバ106はリンク集を利用者に提供し、利用者はWWWブラウザ107を介してリンク集を表示させ、入力指示を行う。

【0140】

これにより、コンテンツの内容にそった質の高いリンク集を簡便に作成することが可能となる。

ところで、図1、26及び27の文書検索装置並びに図28及び29のリンク集生成システムは、図30に示すような情報処理装置（コンピュータ）を用いて構成することができる。図30の情報処理装置は、CPU121、メモリ122、入力装置123、出力装置124、外部記憶装置125、媒体駆動装置126、及びネットワーク接続装置127を備え、それらはバス128により互いに接続されている。

【0141】

メモリ122は、例えば、ROM（Read Only Memory）、RAM（Random Access Memory）等を含み、処理に用いられるプログラムとデータを格納する。CPU121は、メモリ122を利用してプログラムを実行することにより、必要な処理を行う。

【0142】

図1、26及び27の文書検索装置並びに図28及び29のリンク集生成シス

テムの各処理装置を構成する各機器及び各部は、それぞれメモリ 1 2 2 の特定のプログラムコードセグメントにプログラムとして格納される。

【 0 1 4 3 】

入力装置 1 2 3 は、例えば、キーボード、ポインティングデバイス、タッチパネル等であり、ユーザからの指示や情報の入力に用いられる。出力装置 1 2 4 は、例えば、ディスプレイやプリンタ等であり、利用者への問い合わせ、処理結果等の出力に用いられる。

【 0 1 4 4 】

外部記憶装置 1 2 5 は、例えば、磁気ディスク装置、光ディスク装置、光磁気ディスク装置等である。この外部記憶装置 1 2 5 に上述のプログラムとデータを保存しておき、必要に応じて、それらをメモリ 1 2 2 にロードして使用することもできる。

【 0 1 4 5 】

媒体駆動装置 1 2 6 は、可搬記録媒体 1 2 9 を駆動し、その記録内容にアクセスする。可搬記録媒体 1 2 9 としては、メモ리카ード、フロッピーディスク、C D - R O M (C o m p a c t D i s k R e a d O n l y M e m o r y) 、光ディスク、光時期ディスク等、任意のコンピュータ読み取り可能な記録媒体が用いられる。この可搬記録媒体 1 2 9 に上述のプログラムとデータを格納しておき、必要に応じて、それらをメモリ 1 2 2 にロードして使用することもできる。

【 0 1 4 6 】

ネットワーク接続装置 1 2 7 は、LAN (L o c a l A r e a N e t w o r k) 、WAN (W i d e A r e a N e t w o r k) 等の任意のネットワーク (回線) を介して外部の装置を通信し、通信に伴うデータ変換を行う。また、必要に応じて、上述のプログラムとデータを外部の装置から受け取り、それらをメモリ 1 2 2 にロードして使用することもできる。

【 0 1 4 7 】

図 3 1 は、図 3 0 の情報処理装置にプログラムとデータを供給することのできるコンピュータ読み取り可能な記録媒体を示している。可搬記録媒体 1 2 9 や外

部のデータベース 1 3 0 保存されたプログラムとデータは、メモリ 1 2 2 にロードされる。そして、CPU 1 2 1 は、そのデータを用いてそのプログラムを実行し、必要な処理を行う。

【0 1 4 8】

【発明の効果】

本発明によれば、WWWページの重要度を計算する時に、リンク関係だけでなく、URLの類似度を考慮に入れることで、多くのページを抱える特定サイトやそのミラーサイトの重要度が過大に評価されにくくして、重要度を計算することが可能となる。これにより、より精度高く重要文書を選定することが可能となる。

【0 1 4 9】

また、本発明により計算される重要度には悪意を持った個人が意図的に操作しにくい性質がある。

また、本発明によれば、索引として、キーワードの読み（又はスペル）の1文字以上の断片を次々クリックするだけで、効率よくキーワード又は文書データにアクセスすることが可能となる。

【0 1 5 0】

また、キーワード索引において、1画面上のキーワードや文書数を一定個数内に収めることができ、利用者はキーワードを選択しやすいだけでなく、携帯端末のような表示領域が限られたメディアでの利用にも有効である。

【0 1 5 1】

また、本発明によれば、文書データ中のキーワード出現頻度と上述の文書データの重要度に基づいて文書関連度を計算し、この関連度の順に文書データへアクセスするリンクを並べることにより、特定のキーワードに関連する良質の文書への迅速なアクセスを可能とするリンク集を作成することが可能となる。

【0 1 5 2】

本発明によれば、URLの類似と参照／被参照数を利用して、文書のタイプ（メニュー、リンク集、コンテンツ）を判別することが可能となる。さらに、文書タイプの判別結果に基づいて文書データを選定し、選定された文書データについ

てリンク重要度の算出結果及び・又はキーワード出現頻度を組み合わせて、より質の高い文書へのアクセスを可能とするリンク集を作成することが可能となる。

【 0 1 5 3 】

また、本発明によれば、特定のURLを判別することにより、高い精度で特定の分野の文書データを自動的に選定することが可能となる。さらに上述のリンク重要度と判別された特定のURLに基づいて、特定分野の文書データにアクセスすることを可能とするリンク集が精度良く簡単に作成することが可能となる。

【 0 1 5 4 】

またさらに、上述のようにして判別した文書タイプを元に、特定URLの文書データの中から特定の文書タイプの文書データを選択し、これら選択された文書データについての上述のリンク重要に基づいてリンク集を作成することにより、より質の高い特定分野の文書データにアクセスすることを可能とするリンク集を作成することが可能となる。

【図面の簡単な説明】

【図 1】

第 1 実施形態に係わる文書検索装置の構成図である。

【図 2】

文書情報を格納するテーブル集を示す図である。

【図 3】

キーワード情報を格納するテーブル集を示す図である。

【図 4】

索引情報を格納するテーブル集を示す図である。

【図 5】

アクセスログを示す図である。

【図 6】

索引生成処理の手順を示すフローチャートである。

【図 7】

リンク重要度付与器において行う計算を模式的に示す図（その 1）である。

【図 8】

リンク重要度付与器において行う計算を模式的に示す図（その 2）である。

【図 9】

リンク重要度を算出する際に URL 類似度の概念を導入した結果を示す図である。

【図 1 0】

キーワード文字列グラフの一例を示す図である。

【図 1 1】

初期キーワード文字列グラフの生成手順を示すフローチャートである。

【図 1 2】

初期キーワード文字列グラフの生成手順を実現するアルゴリズムの一例を示す図である。

【図 1 3】

中間ノードの縮退処理の手順を示すフローチャートである。

【図 1 4】

中間ノードの縮退処理の手順を実現するアルゴリズムの一例を示す図である。

【図 1 5】

末端ノードの縮退処理の手順を示すフローチャートである。

【図 1 6】

末端ノードの縮退処理の手順を実現するアルゴリズムの一例を示す図である。

【図 1 7】

末端ノードの縮退処理をしたキーワード文字列グラフの一例を示す図である。

【図 1 8】

索引画面の遷移を示す図である。

【図 1 9】

索引トップ画面の一例を示す図である。

【図 2 0】

索引トップ画面のさらなる例を示す図である。

【図 2 1】

索引中間画面の一例を示す図である。

【図 2 2】

索引中間画面のさらなる例を示す図である。

【図 2 3】

索引中間画面のさらなる例を示す図である。

【図 2 4】

キーワード情報画面の一例を示す図である。

【図 2 5】

キーワード情報画面のさらなる例を示す図である。

【図 2 6】

第 2 実施形態に係わる文書検索装置の構成図である。

【図 2 7】

第 3 実施形態に係わる文書検索装置の構成図である。

【図 2 8】

第 4 実施形態に係わるリンク集生成システムの構成図である。

【図 2 9】

第 5 実施形態に係わるリンク集生成システムの構成図である。

【図 3 0】

情報処理装置の構成図である。

【図 3 1】

記録媒体を示す図である。

【図 3 2】

自明な読み入力インターフェースの一例を示す図である。

【符号の説明】

1 1、1 0 2 処理装置

1 2、1 2 3 入力装置

1 3 表示装置

2 1 リンク重要度付与器

2 2 キーワード抽出器

2 3 キーワード文書関連度計算器

- 2 4 索引生成器
- 2 5 索引アクセス部
- 2 6 アクセス解析器
- 2 7 URL類似度計算器
- 3 0 文書データ
- 3 1 リンク重要度
- 3 2 ページキーワード
- 3 3 索引データ
- 3 4 アクセスログ
- 4 1 文書情報テーブル
- 4 2 被参照文書テーブル
- 5 1 キーワードテーブル
- 5 2 キーワード対応テーブル
- 5 3 出現文書テーブル
- 6 1 索引情報テーブル
- 6 2 関連文書テーブル
- 6 3 関連キーワードテーブル
- 7 1 アクセスログ
- 8 1、1 0 1 収集ロボット
- 8 2 同義語データ
- 8 3、1 0 7 WWWブラウザ
- 9 1、1 1 1 文書タイプ判別器
- 9 2、1 1 2 文書タイプ
- 1 0 3 URL文字列判別器
- 1 0 4 特定URL
- 1 0 5 リンク集
- 1 0 6 WWWサーバ
- 1 2 1 CPU
- 1 2 2 メモリ

1 2 4 出力装置

1 2 5 外部記憶装置

1 2 6 媒体駆動装置

1 2 7 バス

1 2 8 ネットワーク接続装置

1 2 9 可搬記録媒体

1 3 0 プログラムデータ

P、P₁、P₂、P₃、q、r₁、r₂、s₁、s₂ ページ

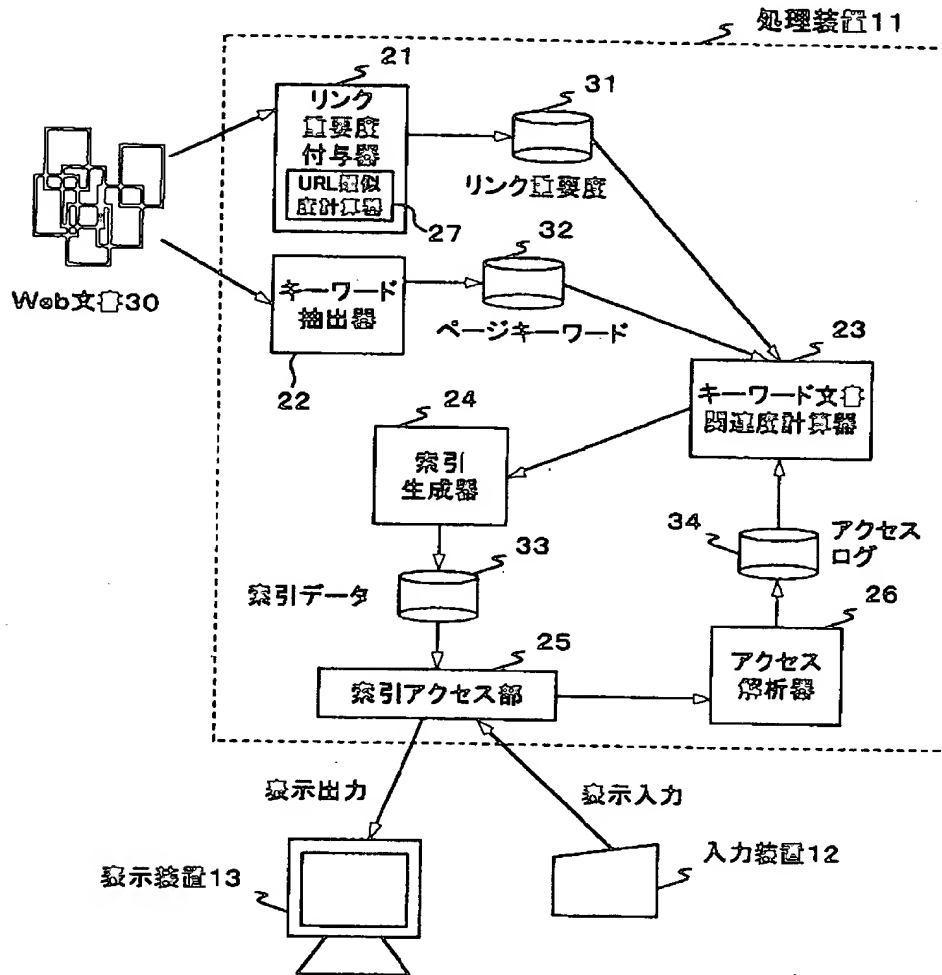
S ステップ

【書類名】

図面

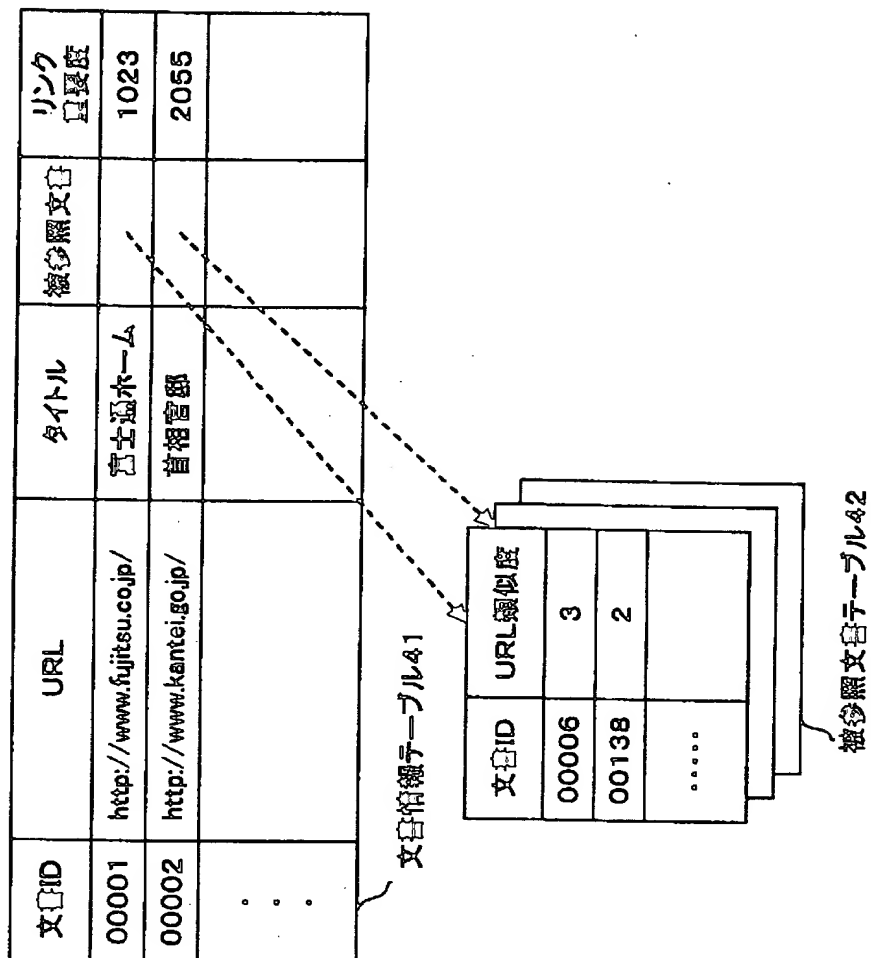
【図 1】

第1実施形態に係わる文書検索装置の構成図



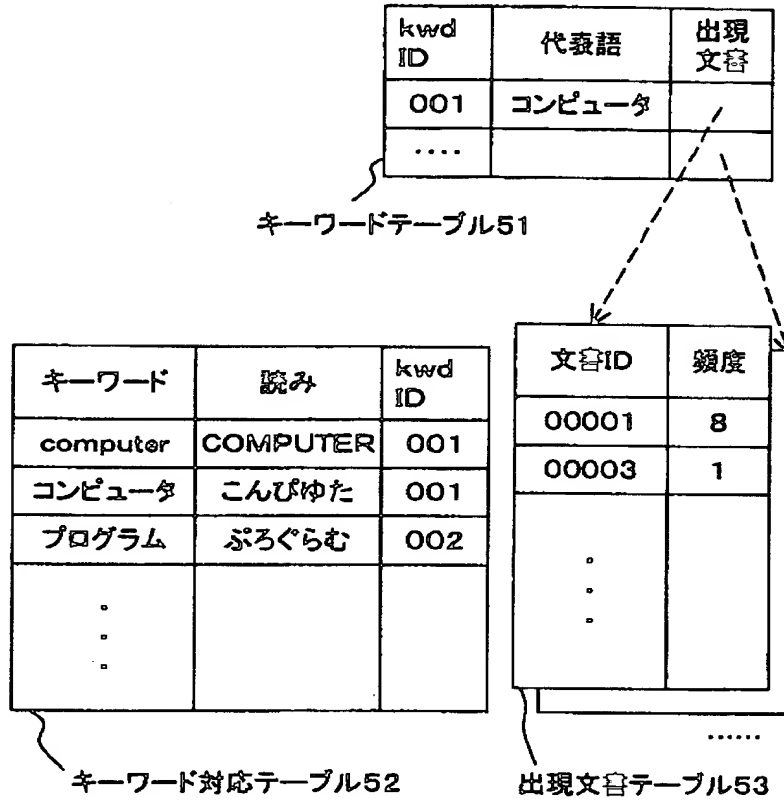
【図 2】

文書情報を格納するテーブル集を示す図



【図 3】

キーワード情報を格納するテーブル集を示す図



【図 4】

索引情報を格納するテーブル集を示す図

文字列	継続する文字列	キーワード列
top	あ, い, ...	
あ	あいぼ, あお, ...	
あいぼ		相棒, アイボリー,
あお	あおぞ	青, 蒼, ...

索引情報テーブル61

キーワード ID	関連文書ID列
093	0005, 0008, ...
321	0004, 0008, ...
....	

関連文書テーブル62

文書ID	関連キーワードID列
0005	093, 099, 122, ...
0008	093, 156, 321, ...
.	
.	
.	

関連キーワードテーブル63

【図5】

アクセスログを示す図

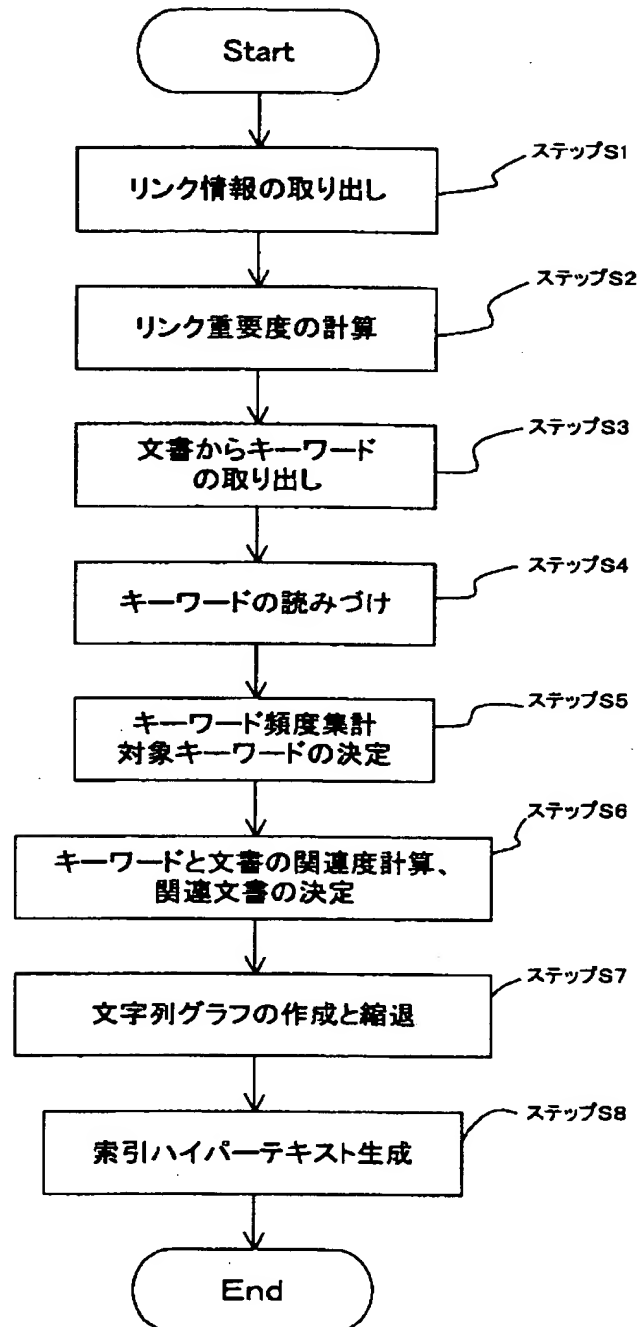
yyyymmddHHMM
の形式

日時	kwd ID	文書ID
200001121436	003	00123
200001121437	005	00054
⋮		
⋮		
⋮		

アクセスログ71

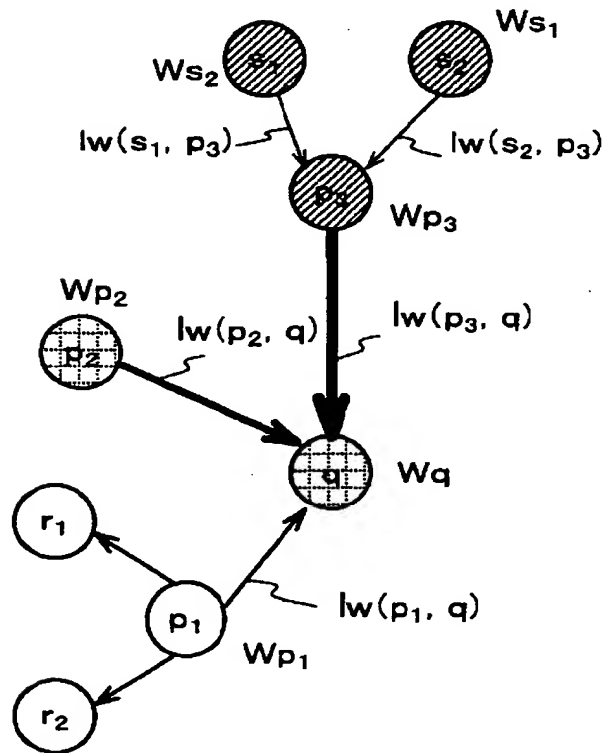
【図 6】

検索処理の手順を示すフローチャート



【図 7】

リンク重要度付与器において行う計算を
模式的に示す図(その1)

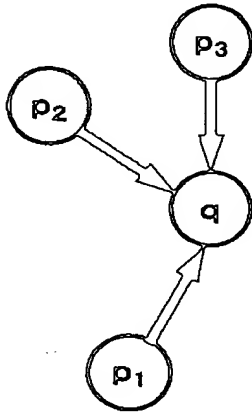


○印:wwwページを示す
○印の太さ:リンクの重みを示す
矢印の太さ:リンクの重みを示す
/////:URLの類似

【図 8】

リンク重要度付与器において行う計算を模式的に示す図(その2)

$$\text{sim}(p_i, q) = 1$$

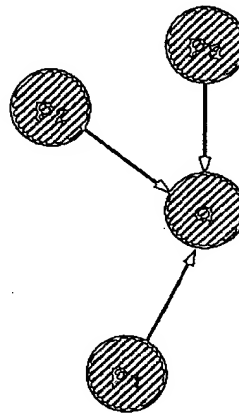


$$lw(p_i, q) = \frac{1}{\text{sim}(p_i, q)} = 1$$

$$w_q = c_q + w_{p1} + w_{p2} + w_{p3}$$

図 8(a)

$$\text{sim}(p_i, q) = n+1$$



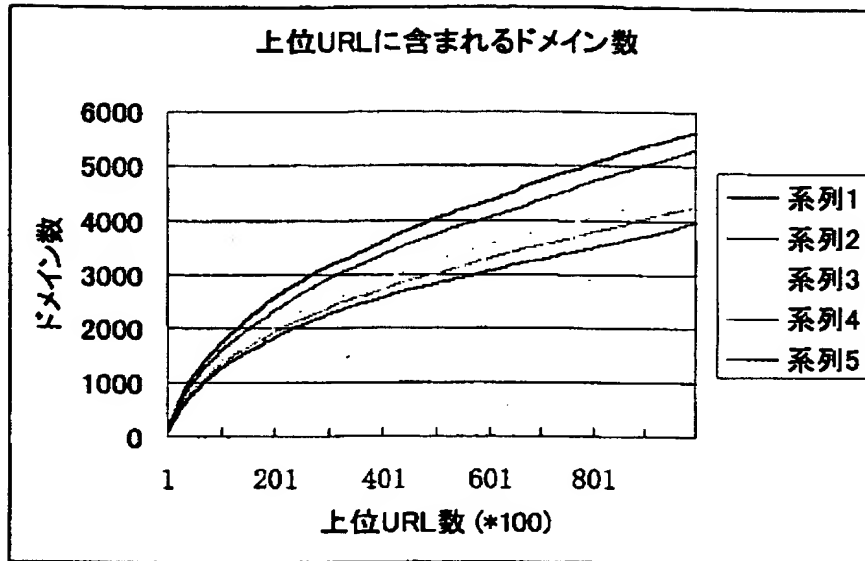
$$lw(p_i, q) = \frac{1}{\text{sim}(p_i, q)} = \frac{1}{n+1}$$

$$w_q = c_q + \frac{w_{p1} + w_{p2} + w_{p3}}{n+1}$$

図 8(b)

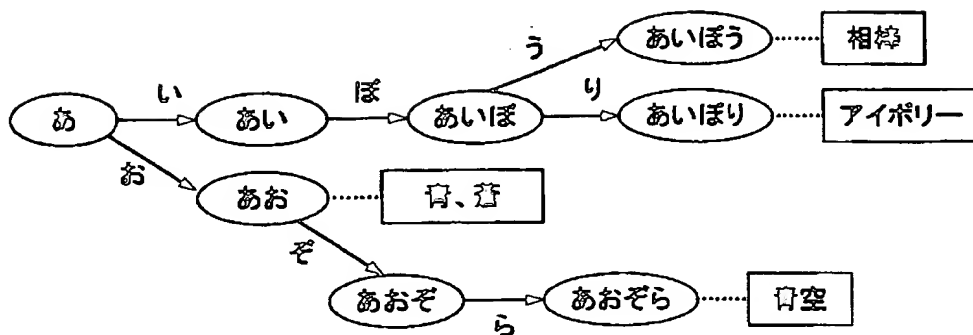
【図 9】

リンク重要度を算出する際に
URL 類似度の概念を導入した結果を示す図

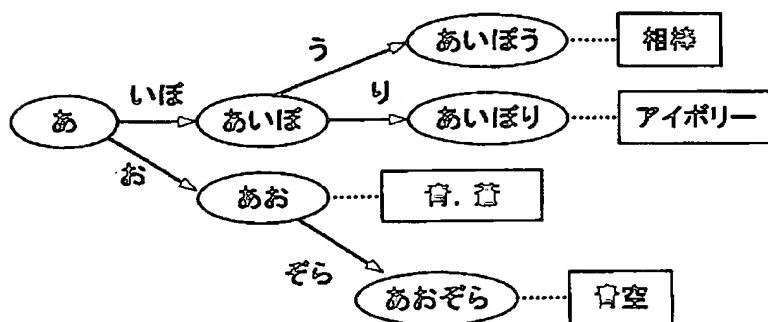


【図 10】

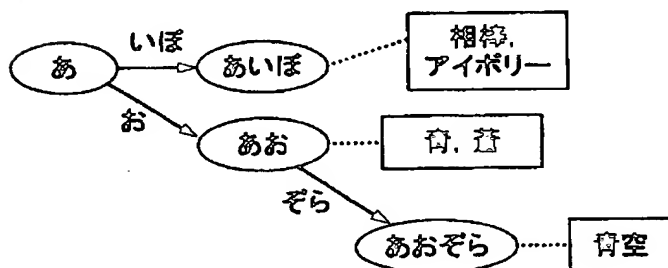
キーワード文字列グラフの一例を示す図



10(a): 初期キーワード文字列グラフ

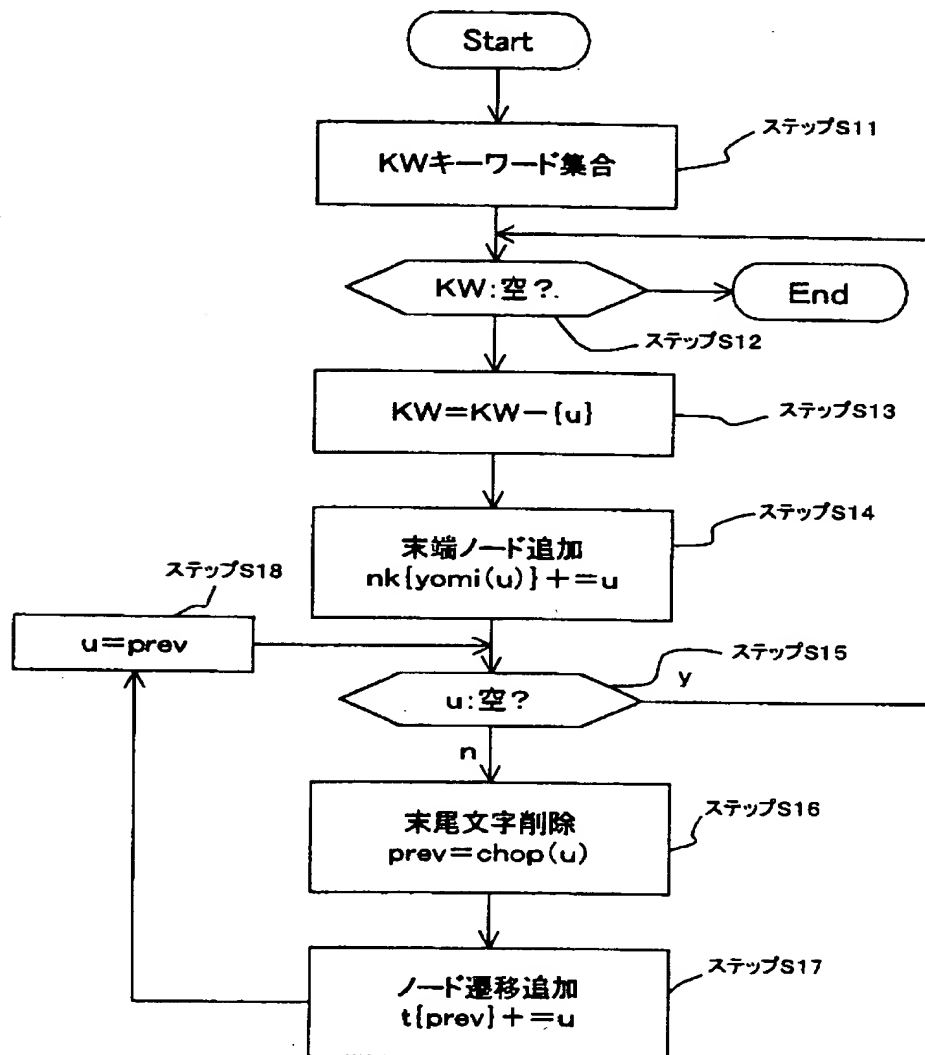


10(b): 中間パスの縮退後のグラフ



10(c): 末端ノードキーワードの縮退後のグラフ

【図 11】

初期キーワード文字列グラフの
生成手順を示すフローチャート

【図 1 2】

初期キーワード文字列グラフ生成手順を
実現するアルゴリズムの一例を示す図

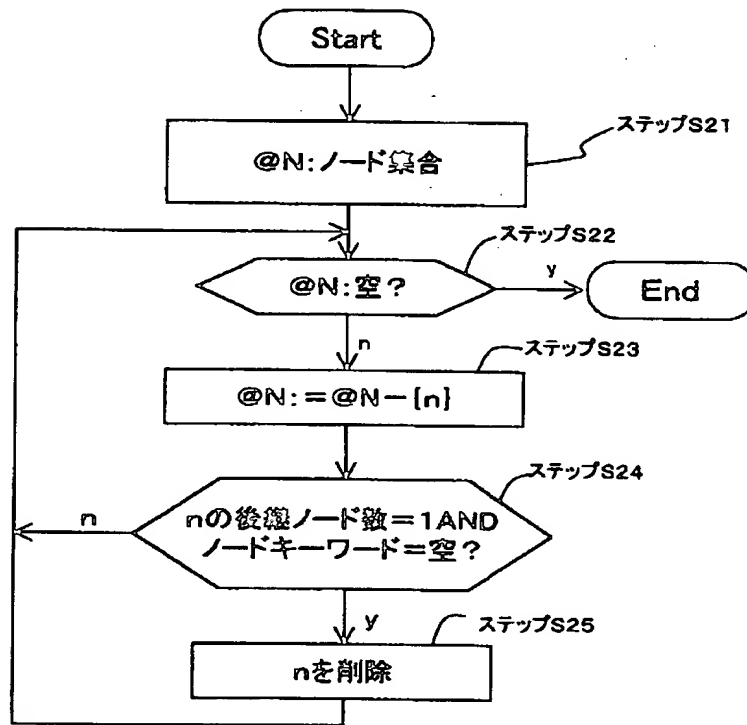
```

proc init_kw_graph ()
{
  @KW:set of keywords;    # キーワード集合
  yomi : YOMI/Spell of keywords; # キーワードの読みを返す関数or 配列
  foreach u in KW {       # 各キーワードuについて
    nk{yomi{u}} . = u. "~+"; # u の読みノードの nk() を設定
    for ( i=0; i<length(u); i++) { # u の文字列長ぶん繰り返し
      local prev = chop(u);        # u の末尾を切取ると親ノード
      t{prev} . = u. "~+";
      u = prev;
    }
  }
}

```

【図 13】

中間ノードの縮退処理の手順を示すフローチャート



【図 14】

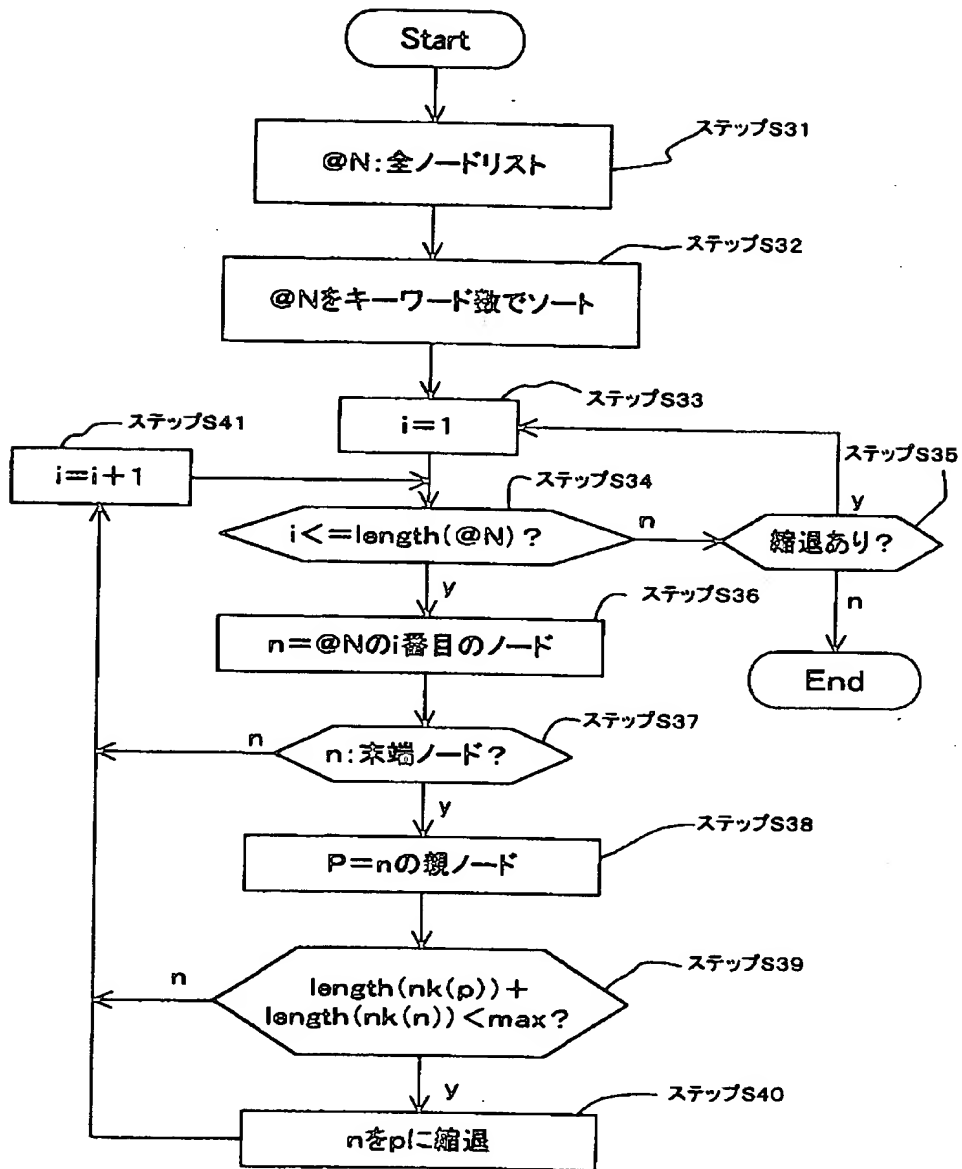
中間ノードの縮退処理の手順を実現する
アルゴリズムの一例を示す図

```

proc shrink_middle ()
{
  @N : set of nodes
  foreach n (@N) {
    next = t{n}: # 次のノードリスト
    kw = nk{n}: # キーワードリスト
    if (length(next) == 1 && kw == "") {
      delete(n) # ノードn削除
    }
  }
}
  
```

【図 15】

端末ノードの縮退処理の手順を示すフローチャート



【図 16】

末端ノードの縮退処理の手順を
実現するアルゴリズムの一例を示す図

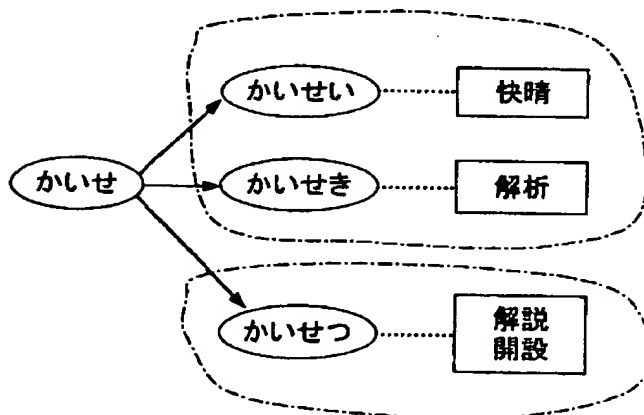
```

proc shrink_leaf ()
{
  @N: set of nodes;    # ノードリスト
  word_max = 2;         # word_max ここでは 2
  changed = true;       # キーワード移動した場合 true
  @N = sort by nk_length @N; # キーワード数の少ない順にソート
  while (changed) {     # 移動がある限り続ける
    changed = false;
    foreach n in @N {
      if (is_leaf(n)) { # 末端ノードの場合
        p = parent_node(n); # 親ノード
        if (length(nk[p]) + length(nk[n]) < word_max) {
          nk[p] = nk[n] . "+" ; # キーワード移動
          delete (n); # 末端ノード削除
          changed = true; # 移動した証
        }
      }
    }
  }
}

```

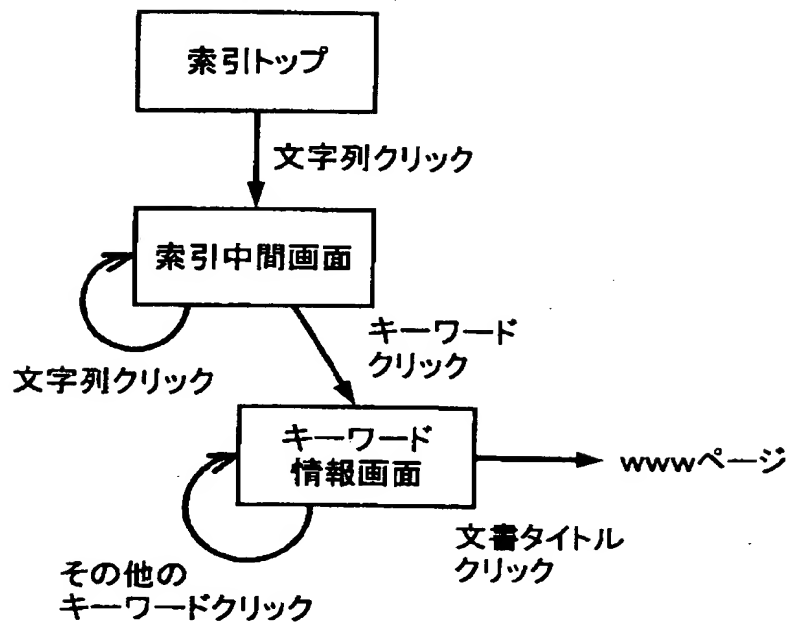
【図 17】

末端ノード縮退処理をした
キーワード文字列グラフの更なる一例を示す図



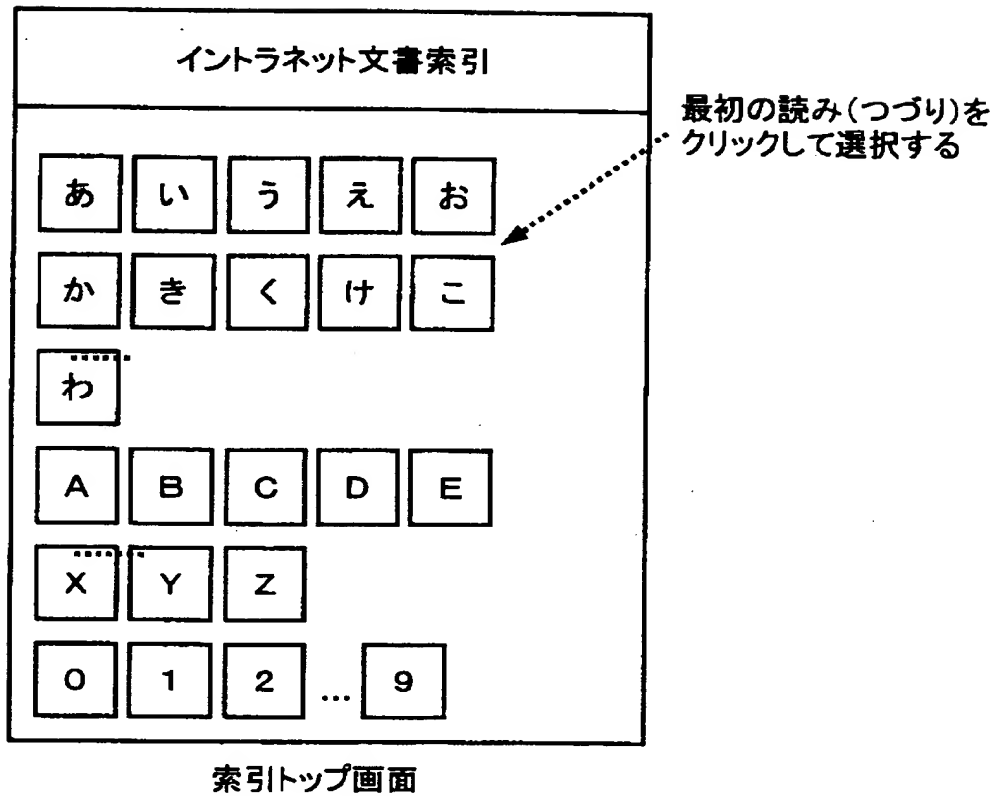
【図18】

索引画面の遷移を示す図



【図 19】

索引トップ画面の一例を示す図



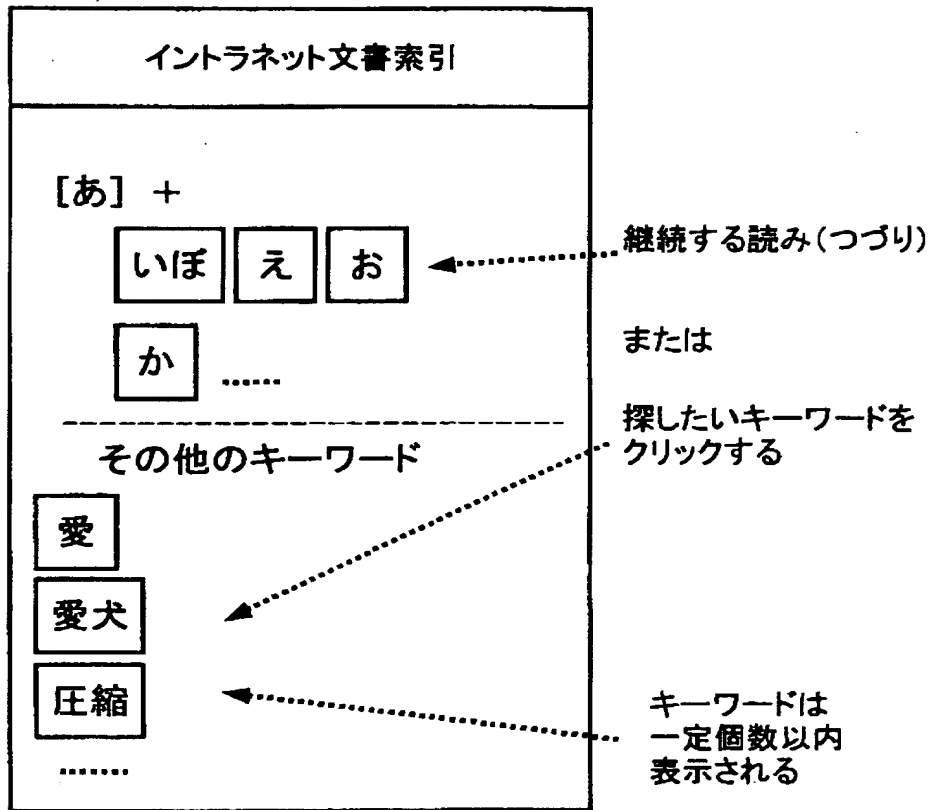
【図20】

索引トップ画面のさらなる例を示す図

社内ページの知的な50音索引		(注意)「ー」(長音)は取って下さい。「っ」「ゃ」「ゆ」を選んで下さい	
あかさたなはまやらわ	1233	あかさたのほちろ	A B C D E F G H I J K L M N O P Q R S T U V W X Y
いまうちはひめゆり	333	うくまつぬふさる	
えはせてねへぬね		きじで	
がきだ	ば	ぐすど	
ばせ	びぶ	こぞ	
はぶ	ひぶ		
ほ	ほ		
び	び		
ぶ	ぶ		
べ	べ		
ぼ	ぼ		

【図 2 1】

索引中間画面の一例を示す図



【图 2 2】

索引中間画面のさらなる例を示す図

「い」

「い」に続く読み

い えろ けこ う が ぎ こ
 か さ た に は め ら ゐ
 ま ち の べ ん と
 り ー じ ゅ ー だ い
 り ゐ ん

その他のキーワード

イオン

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

「い」に続く読み

い えろ けこ う が ぎ こ
 か さ た に は め ら ゐ
 ま ち の べ ん と
 り ー じ ゅ ー だ い
 り ゐ ん

その他のキーワード

イオン

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

イブ

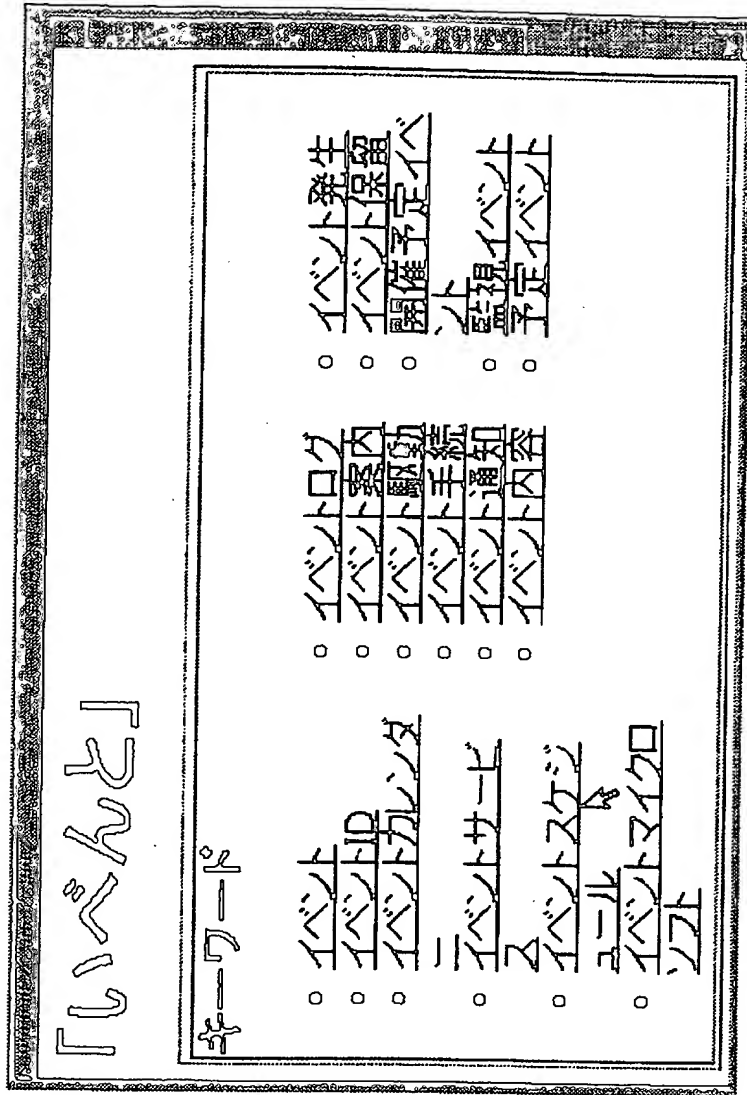
イブ

イブ

イブ

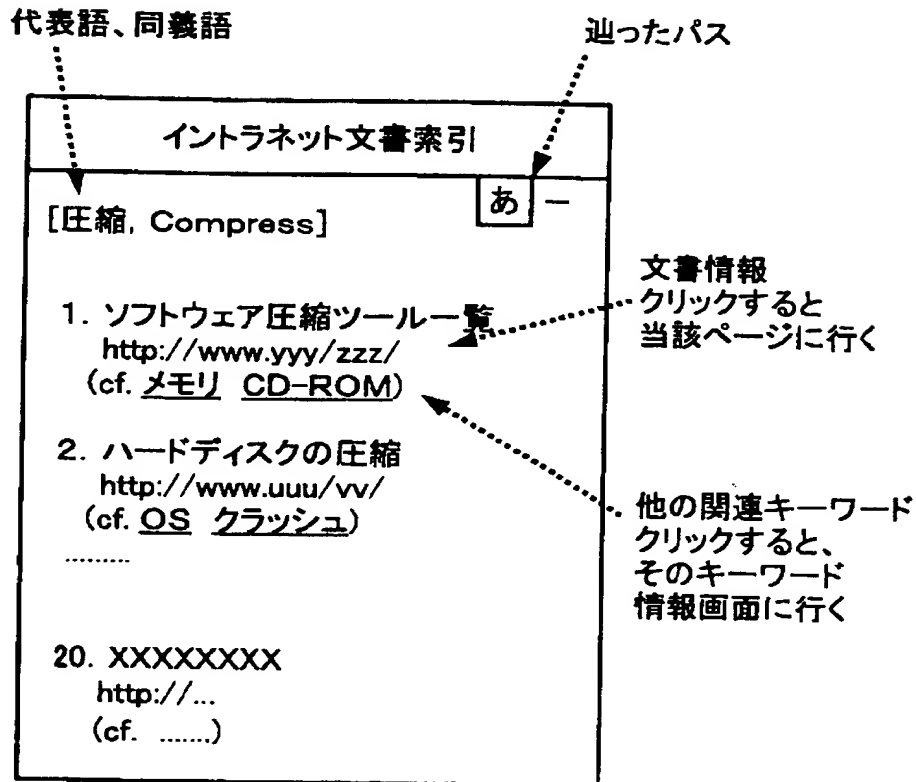
【図23】

索引中間画面のさらなる例を示す図



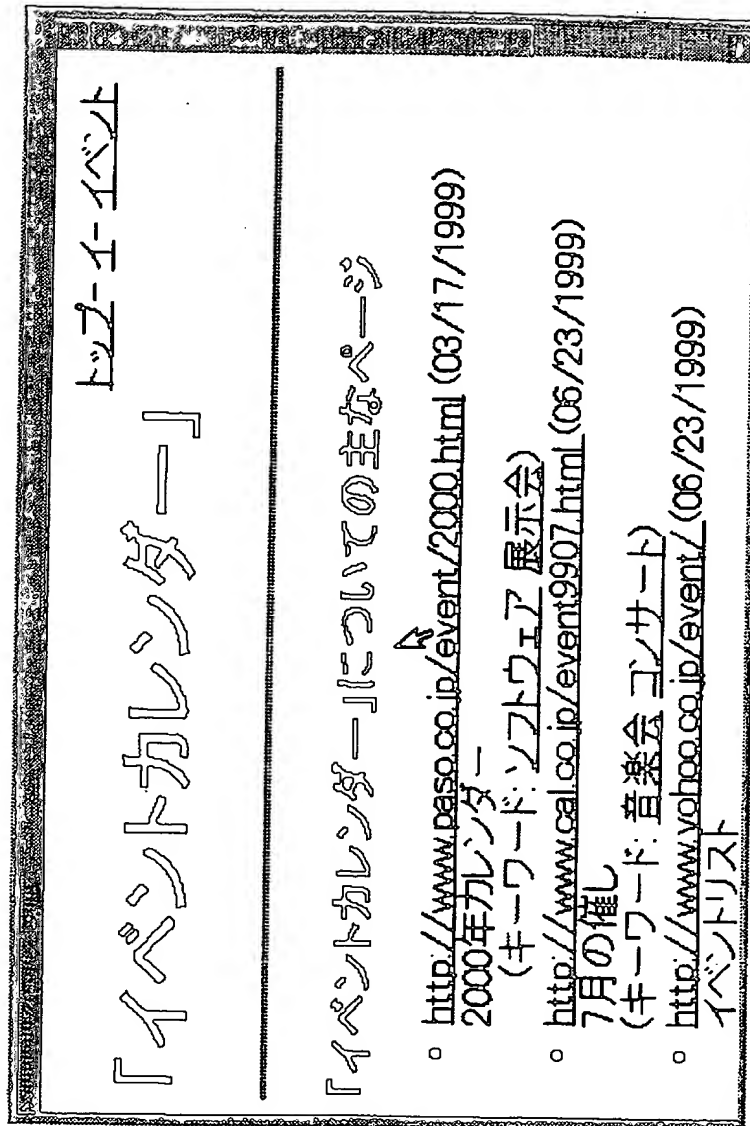
【図 2 4】

キーワード情報画面の一例を示す図



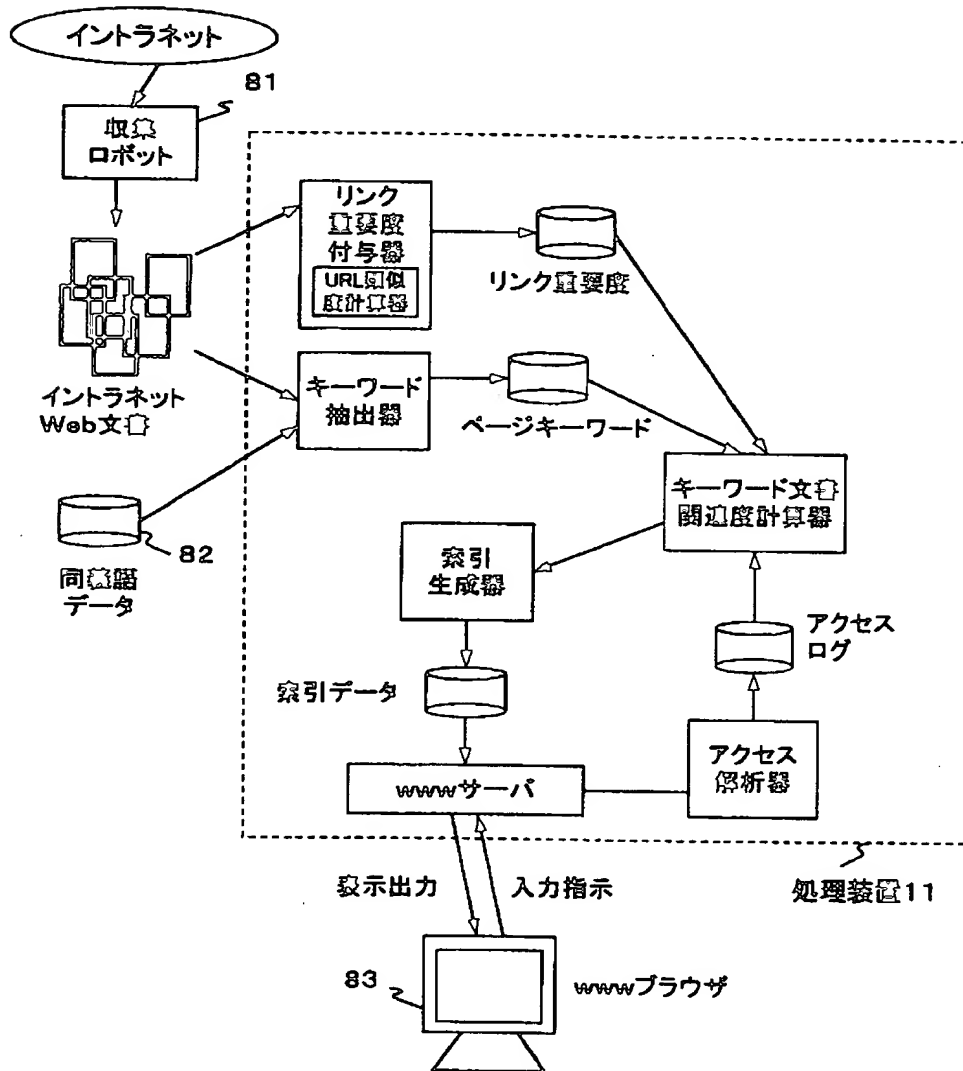
【図25】

キーワード情報画面のさらなる例を示す図



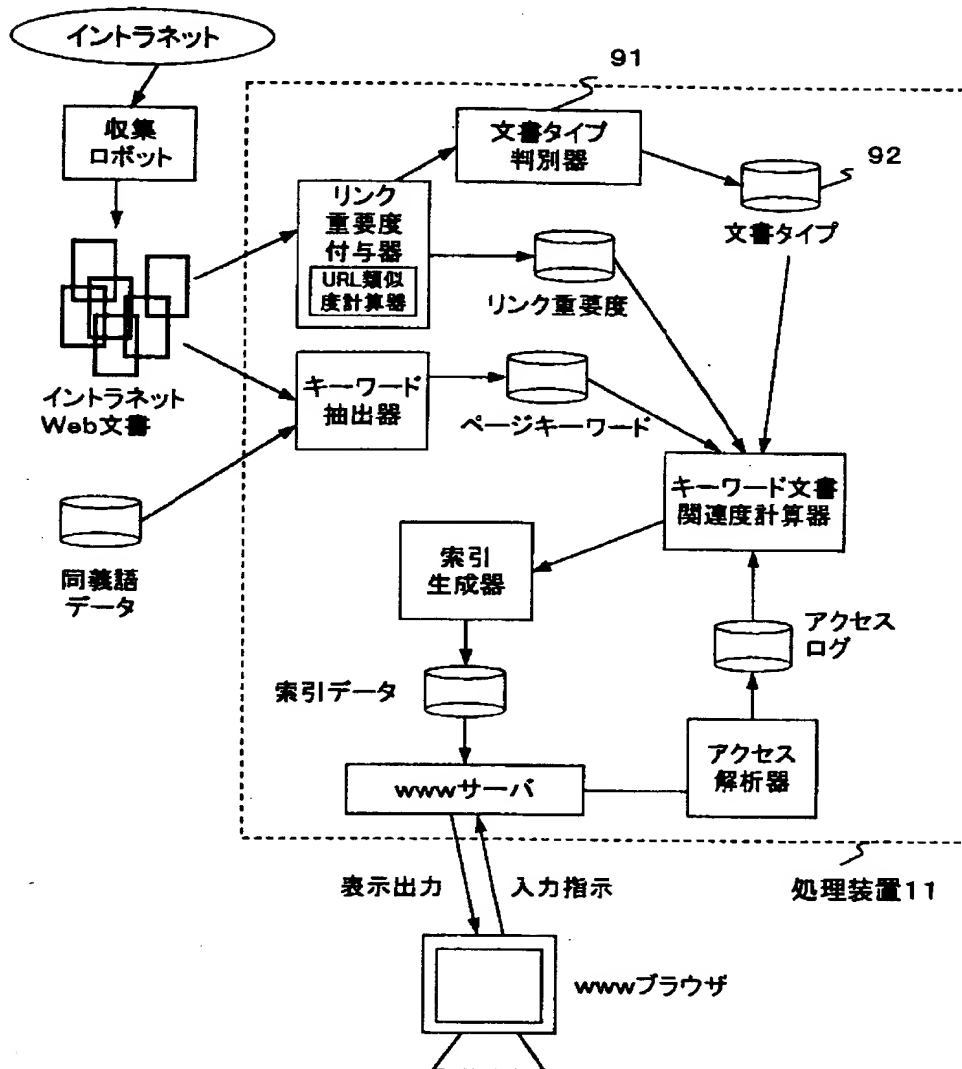
【図 26】

第2実施形態に係わる文書検索装置の構成図



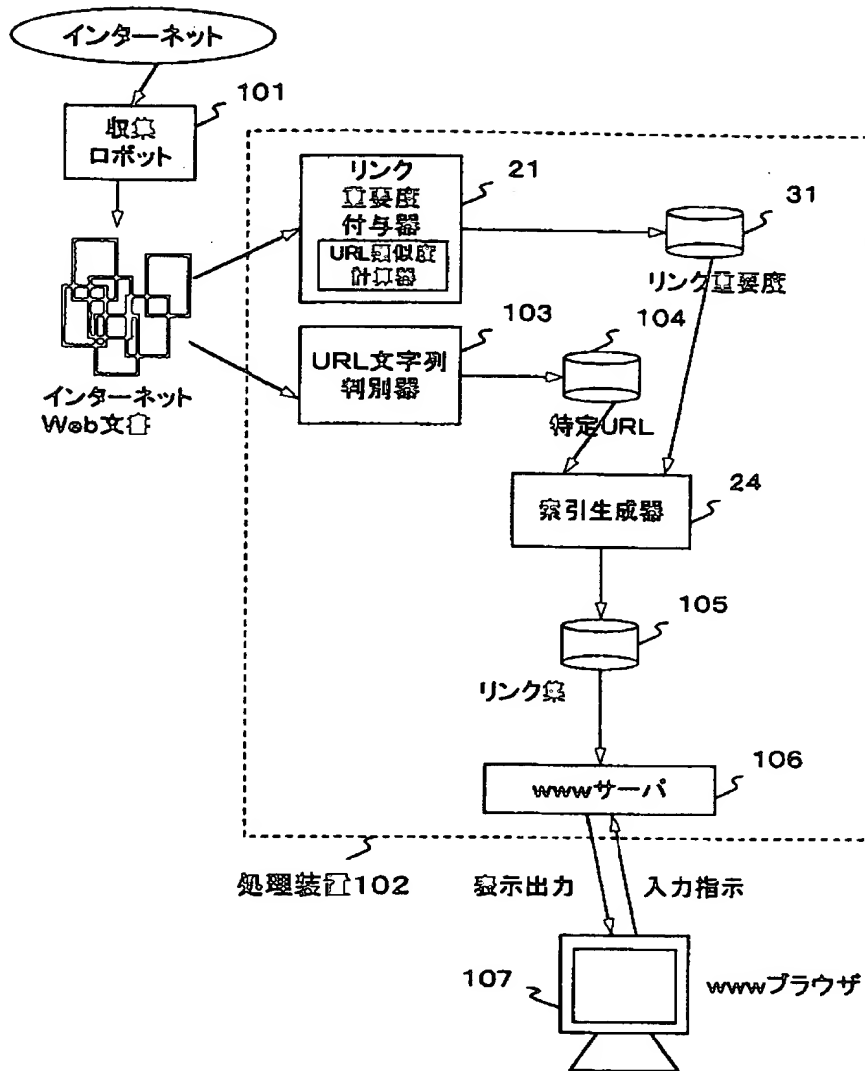
【図27】

第3実施形態に係わる文書検索装置の構成図



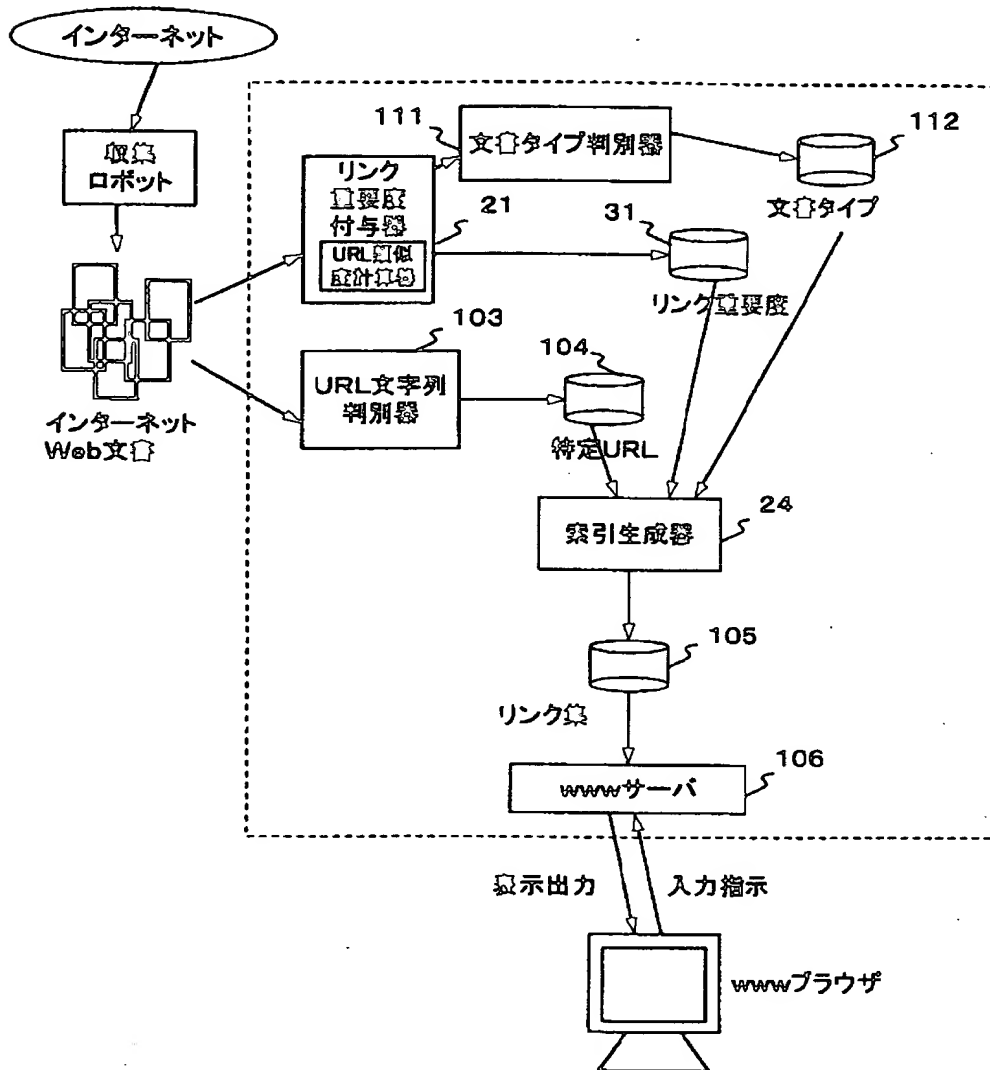
【図 28】

第4実施形態に係わるリンク集生成システムの構成図



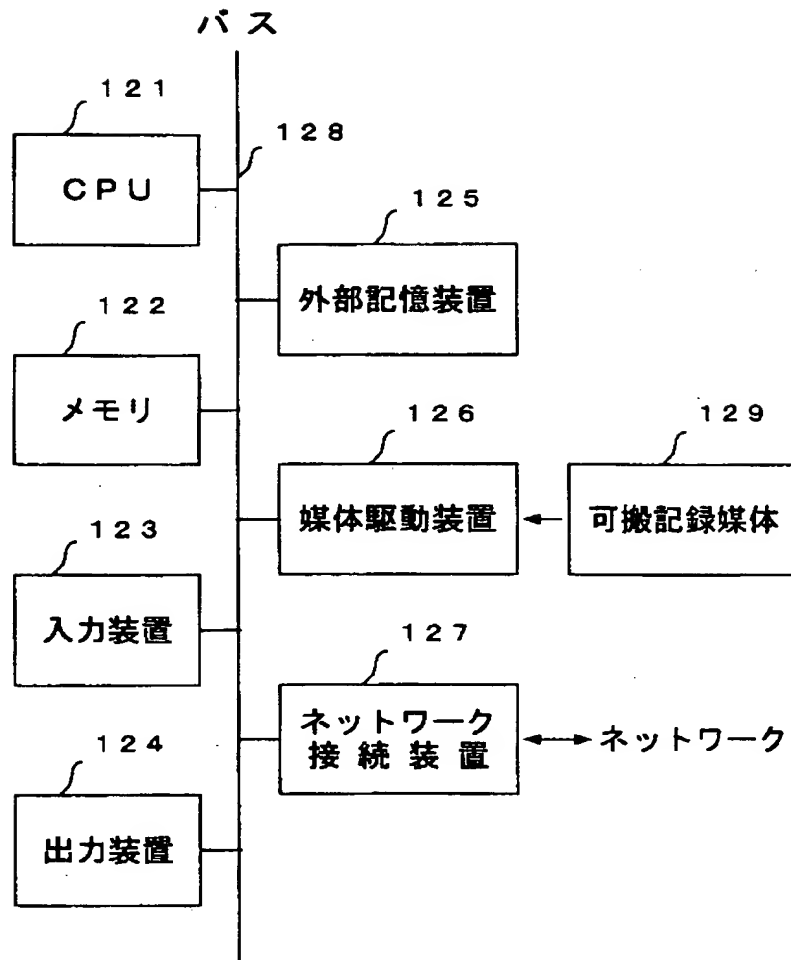
【図 2 9】

第 5 実施形態に係わるリンク集生成システムの構成図



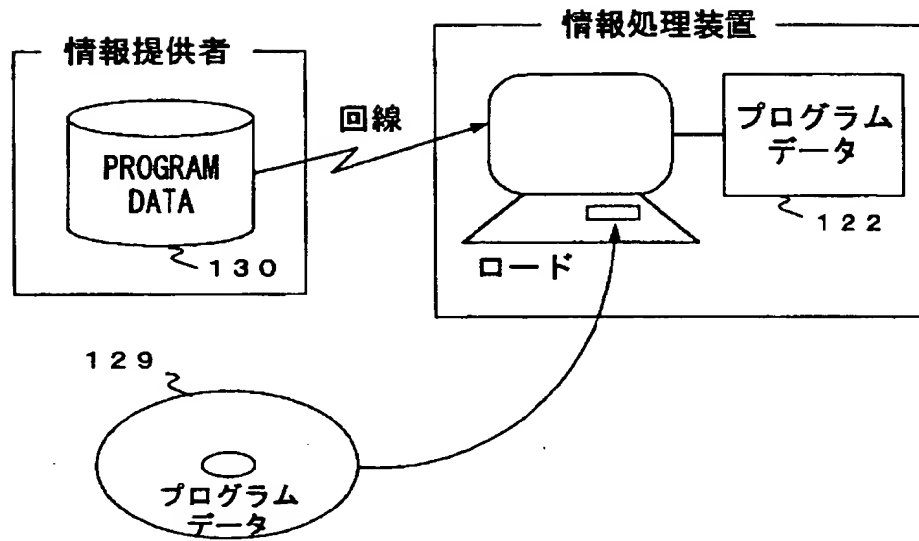
【図 3 0】

情 報 処 理 装 置 の 構 成 図



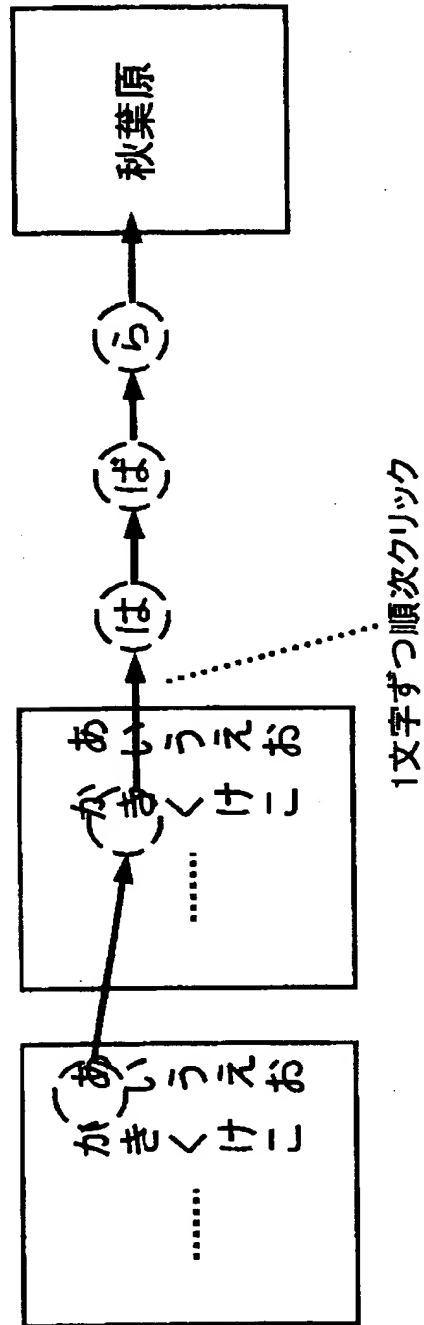
【図 31】

記 録 媒 体 を 示 す 図



【図 3 2】

自明な読み入力インタフェースの
一例を示す図



【書類名】 要約書

【要約】

【課題】 重要な文書データを自動的に検索することを可能とする。

【解決手段】 リンク関係を有する文書データ群から文書データを検索する文書検索装置にリンク重要度付与器 2 1 及びアクセス部 2 5 を備える。リンク重要度付与器 2 1 は文書データ 3 0 のリンク関係に重みを付けて、重要度であるリンク重要度 3 1 を文書データ 3 0 に付与し、アクセス部 2 5 はリンク重要度 3 1 に基づいて文書データ 3 0 にアクセスする。

【選択図】 図 1

出 願 人 履 歴 情 報

識別番号 [000005223]

1. 変更年月日	1996年 3月26日
[変更理由]	住所変更
住 所	神奈川県川崎市中原区上小田中4丁目1番1号
氏 名	富士通株式会社